

Supporting Information for ‘Estimation of Controlled Direct Effects in Longitudinal
Mediation Analyses with Latent Variables in Randomised Studies’

by Wen Wei Loh^{1*}, Beatrijs Moerkerke¹, Tom Loeys¹, Louise Poppe², Geert Crombez²,
and Stijn Vansteelandt^{3,4}

¹ Department of Data Analysis, Ghent University, Ghent, Belgium

² Department of Experimental Clinical and Health Psychology, Ghent University, Ghent,
Belgium

³ Department of Applied Mathematics, Computer Science and Statistics, Ghent
University, Ghent, Belgium

⁴ Department of Medical Statistics, London School of Hygiene and Tropical Medicine,
London, United Kingdom

**email:* WenWei.Loh@UGent.be

Supporting Information for 'Estimation of Controlled Direct Effects in Longitudinal
Mediation Analyses with Latent Variables in Randomised Studies'

Appendix A

G-estimation of the controlled direct effects

In this section, we propose a G-estimator of the controlled direct effects by adapting existing G-estimation methods for time-varying treatments without mediators by Vansteelandt and Sjölander (2016). Using the simplest fitting strategy of Vansteelandt and Sjölander (2016), the G-estimator of the controlled direct effect of X on $Y_t, t = 1, \dots, T$ can be obtained as follows. For $k = t, \dots, 1$:

1. Fit the model for the expected value of M_k , the mediator at time k , by regressing M_k on the history of outcomes \bar{Y}_{k-1} , mediators \bar{M}_{k-1} , confounders \bar{L}_k and treatment X . For example, a linear and additive mediator model for a user-specified link function $g(\cdot)$, allowing for the effects of the predictors on M_k to differ at each time k , is:

$$g\left(\mathbb{E}[M_k|X, \bar{L}_k, \bar{M}_{k-1}, \bar{Y}_{k-1}]\right) = \gamma_{k0} + \gamma_{kx}X + \sum_{j=1}^k \gamma_{kl,j}L_j + \sum_{j=1}^{k-1} \gamma_{km,j}M_j + \sum_{j=1}^{k-1} \gamma_{ky,j}Y_j. \quad (1)$$

For a continuous mediator, let $g(\mu) = \mu$ be the identity link; for a noncontinuous mediator, one may consider other link functions, e.g., the logit link $g(\mu) = \log\{\mu/(1 - \mu)\}$ for a binary mediator.

Denote the regression coefficients by $\gamma_k = (\gamma_{k0}, \dots, \gamma_{ky,k-1})$. Let $\hat{\gamma}_k$ denote the ordinary least squares (OLS) estimators if M_t are continuous, or the maximum likelihood estimators if M_t are noncontinuous. A model that is nonlinear or non-additive, or both, in its predictors is also possible.

Calculate the fitted value of the mediator, denoted by P_k , for each individual by plugging in $\hat{\gamma}_k$ for γ_k in (Eq. 1).

2. If $k = t$, set the outcome, henceforth denoted by R_k , to be $R_k = Y_t$.
3. Fit the linear outcome model by regressing R_k on $\bar{Y}_{k-1}, \bar{M}_{k-1}, \bar{L}_k, X$, as well as P_k and M_k . For example, an additive model that allows for the effects of the

predictors on R_k to differ at each time k is:

$$E[R_k|X, \bar{L}_k, \bar{M}_k, \bar{Y}_{k-1}, P_k] = \beta_{k0} + \beta_{kx}X + \sum_{j=1}^k \beta_{kl,j}L_j + \sum_{j=1}^{k-1} \beta_{km,j}M_j + \sum_{j=1}^{k-1} \beta_{ky,j}Y_j + \beta_{p,k}P_k + \psi_{kt}M_k. \quad (2)$$

An outcome model that is non-additive in its predictors, e.g., by including interaction terms, is also possible. Calculate the OLS estimator of ψ_{kt} , denoted by $\hat{\psi}_{kt}$, and determine the residual $R_{k-1} = R_k - \hat{\psi}_{kt}M_k$.

Repeat steps 1 to 3 until R_0 is obtained. Then fit the regression model

$E[R_0|X] = \alpha_0 + \alpha_1X$. The OLS estimator $\hat{\alpha}_1$ is a G-estimator of the controlled direct effect of X on Y_t . Standard errors can be estimated using a nonparametric percentile bootstrap procedure (Efron & Tibshirani, 1994) that randomly resamples observations with replacement.

Appendix B

Bias-corrected G-estimation method

In this section we develop the G-estimation method for the setting where M_t or Y_t , or both, at each time t are latent but are respectively measured by \tilde{M}_t and \tilde{Y}_t . The G-estimator of the controlled direct effect of X on $Y_t, t = 1, \dots, T$ is obtained in two stages as follows.

Stage 1. Estimate the (Bartlett) scores for the latent variables. Note that when the outcome or the mediator is not latent, then the estimated scores equal the observed values. At each time $t = 1, \dots, T$, for each latent variable:

1. Fit the measurement model in (Eq. 8) to the observed items using *factor analysis* in a structural equation modeling package such as `lavaan` (Rosseel, 2012).

Calculate the maximum likelihood estimates of the measurement model parameters e.g., $\hat{\Lambda}, \hat{\Theta}, \hat{\Sigma}$.

2. Determine the estimated score and estimated measurement error variance by plugging in the estimated parameters for the unknown quantities in (Eq. 9) and in $\Omega = \nu' \Theta \nu$.

Stage 2. Carry out the G-estimation method in Appendix A using the estimated scores in place of the latent variables in each regression model, then apply Fuller's method if any estimated scores are used as predictors. The G-estimator of the controlled direct effect of X on $Y_t, t = 1, \dots, T$ can be obtained as follows. For $k = t, \dots, 1$:

1. Fit the mediator model; e.g., when both the mediator and the outcome are latent,

$$E[\widehat{M}_k | X, \bar{L}_k, \widehat{M}_{k-1}, \widehat{Y}_{k-1}] = \gamma_{k0}^* + \gamma_{kx}^* X + \sum_{j=1}^k \gamma_{kl,j}^* L_j + \sum_{j=1}^{k-1} \gamma_{km,j}^* \widehat{M}_j + \sum_{j=1}^{k-1} \gamma_{ky,j}^* \widehat{Y}_j. \quad (3)$$

The OLS estimator of $\gamma_k^* = (\gamma_{k0}^*, \dots, \gamma_{ky,k-1}^*)'$ is a biased estimator of γ_k in the mediator model where the latent variables had been observed, e.g., in (Eq. 4) with the identity link. To correct the bias using Fuller's method as described in the main text, let W be the latent predictor variable(s) and let Z be the remaining

predictors measured without error in (Eq. 3). For example, if both the mediator and outcome are latent, let $W = (\overline{M}'_{k-1}, \overline{Y}'_{k-1})'$ and $Z = (1, X, \overline{L}'_k)'$. If the outcome is not latent, let $W = \overline{M}_{k-1}$ and $Z = (1, X, \overline{L}'_k, \overline{Y}'_{k-1})'$; if the mediator is not latent, let $W = \overline{Y}_{k-1}$ and $Z = (1, X, \overline{L}'_k, \overline{M}'_{k-1})'$. Denote the resulting bias-corrected estimator using (Eq. 14) by $\hat{\gamma}_k^{MoM}$.

Calculate the fitted value of the mediator, denoted by \hat{P}_k , for each individual by plugging in $\hat{\gamma}_k^{MoM}$ for γ_k^* in (Eq. 3).

2. If $k = t$, set the outcome, henceforth denoted by \hat{R}_k , to be $\hat{R}_k = \hat{Y}_t$.
3. Fit the outcome model to the outcome \hat{R}_k ; e.g., when both the mediator and the outcome are latent,

$$E[\hat{R}_k | X, \overline{L}_k, \widehat{M}_k, \widehat{Y}_{k-1}] = \beta_{k0}^* + \beta_{kx}^* X + \sum_{j=1}^k \beta_{kl,j}^* L_j + \sum_{j=1}^{k-1} \beta_{km,j}^* \widehat{M}_j + \sum_{j=1}^{k-1} \beta_{ky,j}^* \widehat{Y}_j + \beta_{p,k}^* \hat{P}_k + \psi_{kt}^* \widehat{M}_k. \quad (4)$$

When either the mediator or the outcome, or both, are latent, the OLS estimator of ψ_{kt}^* is a biased estimator of ψ_{kt} in the outcome model where the latent variables had been observed, e.g., in (Eq. 5). The bias is corrected using Fuller's method.

The measurement error for the fitted mediator \hat{P}_k can be determined by noting that \hat{P}_k is a linear combination of the estimated scores in the mediator model.

Denote the resulting bias-corrected estimator by $\hat{\psi}_{kt}^{MoM}$. Determine the residual $\hat{R}_{k-1} = \hat{R}_k - \hat{\psi}_{kt}^{MoM} \widehat{M}_k$.

Repeat steps 1 to 3 until $\hat{R}_0 = \hat{Y}_t - \sum_{k=1}^t \hat{\psi}_{kt}^{MoM} \widehat{M}_k$ is obtained. Then fit the regression model $E[\hat{R}_0 | X] = \alpha_0 + \alpha_1 X$. The OLS estimator $\hat{\alpha}_1$ is a G-estimator of the controlled direct effect of X on Y_t . Standard errors can be estimated using a nonparametric bootstrap procedure. Similar to the setting where the mediator and outcome are observed, the fitted mediator \hat{P}_k can be omitted from the outcome model when the outcome model is correctly specified and the predictors in the mediator model are also in the outcome model.

Appendix C

Example R code for path analysis using SEM

The procedure to obtain unbiased estimators of the direct effect of X on Y_t for $t = 1, 2$ is as follows. To illustrate the procedure, example R code used to implement each step is also provided.

1. Fit a saturated path analysis model for all observed variables (i.e., dropping the round nodes and broken lines in the right diagram of Figure 2) using a structural equation modeling package such as `lavaan`. Obtain the maximum likelihood estimates of the edge coefficients.

```

model_obsT2_sem = '
  l1 ~ q1*x
  m1 ~ a1*x + g11*l1
  y1 ~ c1*x + b11*m1 + h11*l1
  l2 ~ q2*x + f12*l1 + k12*m1 + s12*y1
  m2 ~ a2*x + d12*m1 + r12*y1 + g12*l1 + g22*l2
  y2 ~ c2*x + b12*m1 + b22*m2 + e12*y1 + h12*l1 + h22*l2
  # Coefficients for paths for direct effect of X on Y1
  z.CDE_1.gest_sem := c1+q1*h11
  # Coefficients for paths for direct effect of X on Y2
  z.CDE_2.gest_sem := c2+q1*h12+q2*h22+q1*f12*h22+z.CDE_1.gest_sem*(
    e12+s12*h22) '
fit = sem(model=model_obsT2_sem)


```

Note that the set of paths that make up the direct effect of X on Y_t by definition includes paths containing sub-paths for the direct effect of X on prior occurrences $Y_s, s = 1, \dots, t - 1$. E.g., the paths for the direct effect of X on Y_1 consists of the paths $X \rightarrow Y_1$ and $X \rightarrow L_1 \rightarrow Y_1$. The set of paths for the direct effect of X on Y_2 includes the paths $X \rightarrow Y_1 \rightarrow Y_2$ and $X \rightarrow L_1 \rightarrow Y_1 \rightarrow Y_2$ (the direct effect of X on Y_1 followed by the $Y_1 \rightarrow Y_2$ path), as well as $X \rightarrow Y_1 \rightarrow L_2 \rightarrow Y_2$ and $X \rightarrow L_1 \rightarrow Y_1 \rightarrow L_2 \rightarrow Y_2$ (the direct effect of X on Y_1 followed by the $Y_1 \rightarrow L_2 \rightarrow Y_2$ path). It follows that the direct effect of X on Y_2 includes the direct effect of X on Y_1 multiplied by the edge coefficients for $Y_1 \rightarrow Y_2$ (e.g., e_{12}) and for

$$Y_1 \rightarrow L_2 \rightarrow Y_2 \text{ (e.g., } s_{12}, h_{22}\text{)}.$$

2. Combine the estimated path-specific effects for the paths that contribute to the direct effect of X on Y_t for $t = 1, \dots, T$.

```
parest = parameterEstimates( fit )  
parest [grep( "CDE" ,parest$label) ,]
```

The procedure described is equivalent to MLV for $T = 1$.

Appendix D

Number of paths for the direct effects in Figure 2

The number of paths for the direct effect of X on Y_t , denoted by n_t , may be enumerated recursively as follows. There are 2^t possible paths from X to Y_t via \bar{L}_t , including the path with a single $X \rightarrow Y_t$ edge (i.e., without any of \bar{L}_t), that do not intersect any of \bar{Y}_t . To count the number of paths that intersect at least one of \bar{Y}_t , let Y_s be the earliest element in \bar{Y}_t on the path from X to Y_t . Then for $s = 1, \dots, t - 1$, there are 2^{t-s-1} paths from Y_s to Y_t via $\{Y_{s+1}, \dots, Y_{s+(t-s-1)}\}$ (including the path with a single $Y_s \rightarrow Y_t$ edge), and $2^{t-s} - 1$ paths via at least one of $\{L_{s+1}, \dots, L_{s+(t-s)}\}$, for a total of $2^{t-s-1} + (2^{t-s} - 1) = 3(2^{t-s-1}) - 1$ paths. The total number of paths for the direct effect of X on Y_t , $t \geq 2$ can then be calculated recursively as $n_t = 2^t + \sum_{s=1}^{t-1} \{3(2^{t-s-1}) - 1\} n_s$, where $n_1 = 2$. For example, there are $n_2 = 8$ unique direct effect paths for Y_2 ; for $t = 3$, there are $n_3 = 8 + 5n_1 + 2n_2 = 34$ unique direct effect paths.

Appendix E

Simulation studies

Four simulation studies across different settings were conducted to assess the finite-sample biases of the proposed G-estimators empirically. In study 1, the G-estimation method was compared to the path analysis method using a correctly-specified joint model when both the mediator and outcome were observed and continuous. In addition, the path analysis estimators based on a misspecified joint model were calculated. In studies 2 and 3, the G-estimation method was compared to the path analysis method when either the mediator M_t , or the post-treatment confounder L_t , was binary, and both the mediator and outcome were observed. In study 4, the setting where both the mediator and the outcome were latent and continuous was considered. The empirical biases of the estimators using the proposed two-stage G-estimation method were then evaluated under different assumed measurement models for the latent variables.

Study 1

In this study, the mediator and outcome were assumed to be observed and continuous. The G-estimators of the controlled direct effects were compared to two different path analysis estimators: one using the correctly specified (saturated) path model, and another using a path model where the conditional mean models for L_t were misspecified. The simulation study was conducted by carrying out the following steps:

1. A dataset was generated based on the linear path model in the right diagram of Figure 2 for $T = 3$. The conditional mean models are described in Appendix F. Let $X \sim \text{Bernoulli}(0.5)$, $U \sim \mathcal{N}(0, \sigma_u^2)$, $M_0 \sim \mathcal{N}(0, \sigma_m^2)$, $Y_0 \sim \mathcal{N}(0, \sigma_y^2)$, where $\sigma_u^2 = \sigma_m^2 = \sigma_y^2 = 1$. All intercepts were set to zero. The coefficients for the mediator-outcome edges were set to zero (i.e., $b_{11} = b_{12} = b_{13} = b_{22} = b_{23} = b_{33} = 0$), so that all (indirect) effects of treatment on the outcome along paths that intersected any of the mediators were zero. The coefficients for edges emanating from hidden variables U, M_0, Y_0 (i.e., the broken

edges) were set to $-7/9$; the coefficients for all other edges were set to $4/9$. The coefficient values for edges emanating from hidden variables ($-7/9$) were larger in magnitude and had opposite signs as those emanating from observed variables ($4/9$), so that path coefficient estimates might be negative due to hidden confounding. All edge coefficients had absolute magnitude less than one so that the magnitude of each direct effect was smaller than the number of constituent paths (since the effect along any path using the product-of-coefficient method was less than one). The values of $-7/9$ and $4/9$ were chosen merely to reduce the possibility of exact cancellations since their ratio was not an integer (i.e., the former was not an integer multiple of the latter). The errors for all occurrences of the variables were assumed to be independent and normally distributed with mean zero and variance one.

2. The G-estimation procedure was applied to obtain the estimated controlled direct effects.
3. The following saturated path model was fitted to the observed variables using a structural equation modeling package such as `lavaan`. The estimated direct effects were then calculated using the path analysis approach.

```

model_obsT3_sem = '
  l1 ~ q1*x
  m1 ~ a1*x + g11*l1
  y1 ~ c1*x + b11*m1 + h11*l1
  l2 ~ q2*x + f12*l1 + k12*m1 + s12*y1
  m2 ~ a2*x + d12*m1 + r12*y1 + g12*l1 + g22*l2
  y2 ~ c2*x + b12*m1 + b22*m2 + e12*y1 + h12*l1 + h22*l2
  l3 ~ q3*x + f13*l1 + f23*l2 + k13*m1 + k23*m2 + s13*y1 + s23*y2
  m3 ~ a3*x + d13*m1 + d23*m2 + r13*y1 + r23*y2 + g13*l1 + g23*l2 +
    g33*l3
  y3 ~ c3*x + b13*m1 + b23*m2 + b33*m3 + e13*y1 + e23*y2 + h13*l1 +
    h23*l2 + h33*l3
  z.CDE_1.gest_sem := c1+q1*h11
  z.CDE_2.gest_sem := c2+q1*h12+q2*h22+q1*f12*h22+z.CDE_1.gest_sem*(

```

```

e12+s12*h22)
z.CDE_3.gest_sem := c3+q1*h13+q2*h23+q3*h33+q1*f12*h23+q1*f13*h33+q2
*f23*h33+q1*f12*f23*h33+z.CDE_1.gest_sem*(e13+s12*h23+s13*h33+s12
*f23*h33)+z.CDE_2.gest_sem*(e23+s23*h33)
,

```

4. A misspecified path model where the $M_1 \rightarrow L_2$, $M_1 \rightarrow L_3$, $M_2 \rightarrow L_3$ edges were omitted from the saturated path model was then fitted to the observed variables. Note that the constituent paths for the direct effects did not include these edges by definition since these were edges emanating from the mediator. The edge coefficients were constrained to zero in the above path model using `lavaan` as follows:

```

model_constraints = '
# fix edge coefficients for Ms -> Lt where s<t to be zero
k12 == 0
k13 == 0
k23 == 0
,
model_obsT3_sem_restricted = paste(model_obsT3_sem, model_constraints)

```

The estimated direct effects under the misspecified model were then calculated using the path analysis approach.

Steps 1 to 4 were repeated 50000 times each for sample sizes $n = 60, 150, 300$. About 12% of the simulated datasets for $n = 60$ (10% for $n = 150, 300$) where SEM failed to converge and the estimated standard errors for the edge coefficients could not be computed were discarded. The average biases of the three considered estimators for the different samples sizes are plotted in Figure E1. The G-estimators and the path analysis estimators (using the correctly-specified saturated path model) of the controlled direct effects at each time $t = 1, 2, 3$, were identical and unbiased empirically at all sample sizes. However, using a misspecified path model yielded empirically biased path analysis estimators of the direct effects of X on Y_2 and Y_3 , even at larger sample sizes. The misspecified conditional mean models for L_2 and L_3 (due to omitting prior occurrences

of the mediator as predictors) resulted in biased product-of-coefficient (PC) estimates of effects along paths that intersected L_2 and/or L_3 , which contributed to the biases in the direct effects on Y_2 and Y_3 .

The biases of the path analysis estimator (under the misspecified model) that increased in magnitude with t , as observed in Figure E1, can be explained as follows. As previously described in Appendix C, the path analysis estimator of the direct effect of X on $Y_t, t > 1$, by definition always includes the direct effect of X on a previous occurrence of the outcome $Y_s, 0 < s < t$. The direct effect is then multiplied by the effect(s) along path(s) from Y_s to Y_t using the PC method. Thus any biases in the estimated direct effect of X on Y_s are accumulated in the bias of the estimated direct effect of X on Y_t when the effects of Y_s on Y_t are positive.

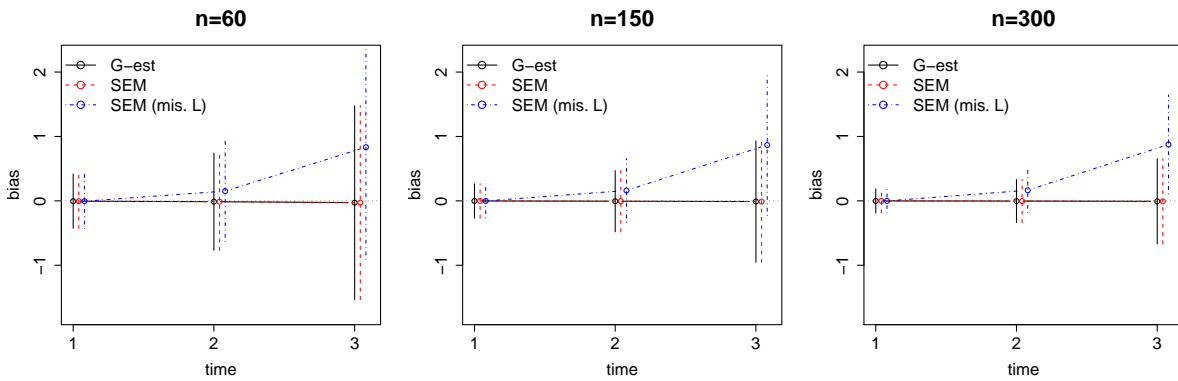


Figure E1. Average biases of the G-estimators ('G-est') and the path analysis estimators, using either a correctly specified path model ('SEM') or a misspecified path model ('SEM (mis. L)'), of the controlled direct effects at each time $t = 1, 2, 3$. Both the mediator and outcome were continuous variables that were observed at all times. The true values of the direct effects were 0.64 ($t = 1$), 1.34 ($t = 2$), 2.80 ($t = 3$). The empirical standard errors are drawn as vertical lines. The path analysis estimators are shifted slightly to the right for clarity. The sample sizes n are stated in each panel.

Study 2

Study 1 was repeated, but with a binary mediator as follows. In the data-generating process in step 1, M_t^* was first randomly sampled from a Bernoulli distribution with probability of success $\Phi(E[M_t|X, \bar{L}_t, \bar{M}_{t-1}, \bar{Y}_{t-1}, U_m])$, where $E[M_t|X, \bar{L}_t, \bar{M}_{t-1}, \bar{Y}_{t-1}, U_m]$ was the conditional mean mediator model described in Appendix F, and $\Phi(\cdot)$ was the cumulative distribution function of a standard normal distribution. Then M_t was set to M_t^* . The path analysis estimator was obtained using SEM where the mediators \bar{M}_3 were specified as noncontinuous ('categorical') variables in `lavaan`, and the maximum likelihood estimates obtained using the three-stage weighted least squares (3SWLS) approach (Browne, 1984):

```
fit = sem(model=model_obsT3_sem, ordered=c("m1", "m2", "m3"))
```

About 6% of the simulated datasets for $n = 60$ (5% for $n = 150, 300$) where SEM failed to converge and the estimated standard errors for the edge coefficients could not be computed were discarded. A smaller proportion of datasets were discarded when M_t was binary than when M_t was continuous (in study 1) as robust standard errors were computed under the former (by using the full weight matrix to correct the standard errors), whereas maximum likelihood estimation for the standard errors was employed in the latter. The average biases and empirical standard errors of the G-estimators and the path analysis estimators (for the datasets where SEM converged successfully) are displayed in Table E1. The path analysis estimators of the controlled direct effects of X on Y_2 and Y_3 were empirically biased, even in large sample sizes, whereas the G-estimators were empirically unbiased (up to Monte Carlo sampling error).

The biases of the path analysis estimators of the direct effects can be explained as follows. When a variable is specified as a binary (or 'categorical') variable in `lavaan`, its underlying continuous measure is used as the response variable. However, if the variable is also an 'endogenous' variable, i.e., it is simultaneously a predictor in another linear model, then the underlying continuous measure is also used (instead of the observed binary measure) as a predictor when fitting the joint model using the 3SWLS approach.

Sample size	Method	Y_1		Y_2		Y_3	
		Bias	E.s.e.	Bias	E.s.e.	Bias	E.s.e.
60	G-estimation	0.00	0.42	-0.01	0.75	-0.02	1.47
	SEM	0.00	0.43	-0.03	0.76	-0.05	1.49
150	G-estimation	0.00	0.26	0.00	0.47	0.00	0.92
	SEM	0.00	0.27	-0.02	0.47	-0.03	0.92
300	G-estimation	0.00	0.19	0.00	0.33	0.00	0.65
	SEM	0.00	0.19	-0.02	0.33	-0.04	0.65

Table E1

Average biases and empirical standard errors (E.s.e.) of estimators of the controlled direct effects of X on $Y_t, t = 1, 2, 3$, when the mediator was binary. The true values of the direct effects were 0.64 ($t = 1$), 1.34 ($t = 2$), 2.80 ($t = 3$).

Maximum likelihood estimates of the regression (path) coefficients in models that include occurrences of the binary mediator $M_t, t = 1, 2, 3$, as predictors may thus be biased (even in large samples). Note that the path analysis estimator of the direct effect on the outcome at $t = 1$ was biased (but only up to the fourth decimal place), since the model for Y_1 included binary M_1 as a predictor.

Study 3

At the request of a reviewer, study 1 was repeated with a binary confounder L_t as follows. The coefficients for all confounder-outcome edges were set to zero (i.e., $h_{11} = h_{12} = h_{13} = h_{22} = h_{23} = h_{33} = 0$), so that the effects along any path from treatment to one of the outcomes that intersected any of (L_1, L_2, L_3) were zero. In other words, the direct effect of X on $Y_t, t = 1, 2, 3$, consisted only of paths that intersected any of the prior occurrences of the outcome \bar{Y}_{t-1} . The binary confounder L_t was similarly generated as the binary mediator in study 2. In the data-generating process in step 1, L_t^* was first randomly sampled from a Bernoulli distribution with probability of success $\Phi(\mathbb{E}[L_t|X, \bar{L}_{t-1}, \bar{M}_{t-1}, \bar{Y}_{t-1}, U])$, where $\mathbb{E}[L_t|X, \bar{L}_{t-1}, \bar{M}_{t-1}, \bar{Y}_{t-1}, U]$ was the

conditional mean model for L_t described in Appendix F, then L_t was set to L_t^* . All occurrences of the mediator and outcome were continuous. The path analysis estimator was obtained using SEM where the variables \bar{L}_3 were specified as noncontinuous ('categorical') variables in `lavaan`, and the maximum likelihood estimates obtained using the 3SWLS approach:

```
fit = sem(model=model_obsT3_sem, ordered=c("11", "12", "13"))
```

About 5% of the simulated datasets where SEM failed to converge and the estimated standard errors for the edge coefficients could not be computed were discarded. The average biases and empirical standard errors of the G-estimators and the path analysis estimators (for the datasets where SEM converged successfully) are displayed in Table E2. Similar to the setting in study 2 with a binary mediator, the path analysis estimators of the controlled direct effects of X on Y_2 and Y_3 were empirically biased, even in large sample sizes, whereas the G-estimators were empirically unbiased (up to Monte Carlo sampling error).

Sample size	Method	Y_1		Y_2		Y_3	
		Bias	E.s.e.	Bias	E.s.e.	Bias	E.s.e.
60	G-estimation	0.00	0.36	0.00	0.47	-0.01	0.65
	SEM	0.02	0.36	0.03	0.48	0.05	0.67
150	G-estimation	0.00	0.22	0.00	0.30	0.00	0.40
	SEM	0.02	0.23	0.03	0.30	0.05	0.41
300	G-estimation	0.00	0.16	0.00	0.21	0.00	0.28
	SEM	0.02	0.16	0.03	0.21	0.05	0.29

Table E2

Average biases and empirical standard errors (E.s.e.) of estimators of the controlled direct effects of X on Y_t , $t = 1, 2, 3$, when the post-treatment confounder L_t was binary. The true values of the direct effects were 0.44 ($t = 1$), 0.64 ($t = 2$), 0.93 ($t = 3$).

Study 4

In this study, the setting where both the mediator and outcome were continuous and latent was considered. All occurrences of the variables (L_t, M_t, Y_t) were first generated for $t = 1, 2, 3$, using the same procedure as step 1 of study 1. The (manifest) items \tilde{M}_t and \tilde{Y}_t were then generated from M_t and Y_t respectively, using specified values of the factor loadings Λ_{mt} and Λ_{yt} , as:

$$\tilde{M}_t = \Lambda_{mt}M_t + \tilde{U}_{mt}, \quad \tilde{Y}_t = \Lambda_{yt}Y_t + \tilde{U}_{yt}; \quad t = 1, 2, 3.$$

There were four items in \tilde{M}_t and in \tilde{Y}_t at each time t . The values of the factor loadings were set to $\Lambda_{mt} = \Lambda_{yt} = (1, 3/4, 3/2, -1)$, which did not depend on the measurement occasion t . The values were chosen merely to induce variability in the factor loadings for different items, since one was positive and smaller than one, another was larger than one, and a third was negative. The measurement errors \tilde{U}_{mt} and \tilde{U}_{yt} were assumed to be independent and normally distributed with mean vector zero and covariance matrix $\Theta = \theta\mathbf{I}_4$ where $\theta = 10/4$ and \mathbf{I}_4 was the 4×4 identity matrix. The value of θ was chosen so that the (conditional) variance of each item was larger than the (conditional) variance of the latent variable (either σ_m^2 or σ_y^2). The items were (conditionally) independent of one another for each occurrence of the latent variable, and across different occurrences of the latent variable. The variables M_t and Y_t were then deleted from each observed dataset after the items \tilde{M}_t and \tilde{Y}_t were generated.

The two-stage G-estimation procedure was applied to obtain the estimated controlled direct effects. In the first stage, the correctly-specified measurement model was fitted to the manifest items using `lavaan` as follows:

```
o_model_v = '
# measurement models for latent mediator at each time t
m1 =~ o_m1_1 + o_m1_2 + o_m1_3 + o_m1_4
m2 =~ o_m2_1 + o_m2_2 + o_m2_3 + o_m2_4
m3 =~ o_m3_1 + o_m3_2 + o_m3_3 + o_m3_4
# measurement models for latent outcome at each time t
y1 =~ o_y1_1 + o_y1_2 + o_y1_3 + o_y1_4
```

```

y2 =~ o_y2_1 + o_y2_2 + o_y2_3 + o_y2_4
y3 =~ o_y3_1 + o_y3_2 + o_y3_3 + o_y3_4
,
fit = cfa(model=o_model_v, orthogonal=TRUE)

```

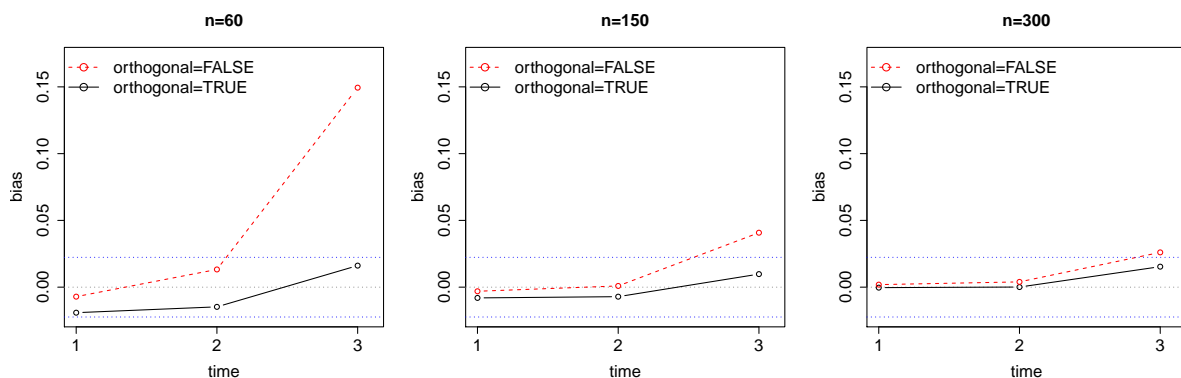


Figure E2. Average biases of the two-stage G-estimators, either assuming an orthogonal covariance matrix for the latent variables (`orthogonal=TRUE`), or allowing for correlated latent variables (`orthogonal=FALSE`), in the fitted measurement model. The true factor loadings were fixed at the same value across occasions. Both the mediator and outcome were continuous variables that were latent at all occasions. The true values of the direct effects were 0.64 ($t = 1$), 1.34 ($t = 2$), 2.80 ($t = 3$). The Monte Carlo sampling errors are drawn as horizontal blue dotted lines. The sample sizes n are stated in each panel.

The average biases of the G-estimators of the controlled direct effects at each time $t = 1, 2, 3$, for the different samples sizes are plotted in Figure E2. The G-estimators were unbiased empirically at all sample sizes (up to Monte Carlo sampling error). The argument `orthogonal=TRUE` constrains all covariances of the latent variables to be zero when fitting the measurement model in the first stage. For comparison, a (misspecified) measurement model where the latent variables were allowed to be correlated was also fitted, by setting the argument `orthogonal=FALSE` (the default setting in `lavaan`) as follows:

```

fit = cfa(model=o_model_v, orthogonal=FALSE)

```

As shown in Figure E2, the G-estimators of the direct effect on Y_3 using such a misspecified measurement model in the first stage were empirically biased at smaller sample sizes, but the biases diminished at larger sample sizes. The observed biases at smaller sample sizes were due to biased estimates of the measurement model parameters in the first stage, when the model with a larger number of free parameters, specifically the covariances between the latent variables, was fitted.

The study was then repeated using factor loadings that depended on the measurement occasion t to generate the items: $\Lambda_{mt} = \Lambda_{yt} = (1, t/4, t/2, -3/t)$. Again the values were chosen merely to induce variability in the factor loadings for different items, since at least one was positive and (strictly) smaller than one, and another was negative and no larger than minus one. The factor loadings at $t = 3$ equalled those in the previous setting. The average biases of the G-estimators assuming a measurement model that was either correct ('orthogonal=TRUE') or incorrect ('orthogonal=FALSE') are plotted in Figure E3. The results were similar to those when the factor loadings did not vary over time: when the measurement model was correctly specified, the G-estimators were unbiased empirically at all sample sizes (up to Monte Carlo sampling error); otherwise the biases due to a misspecified measurement model diminished as the sample size increased.

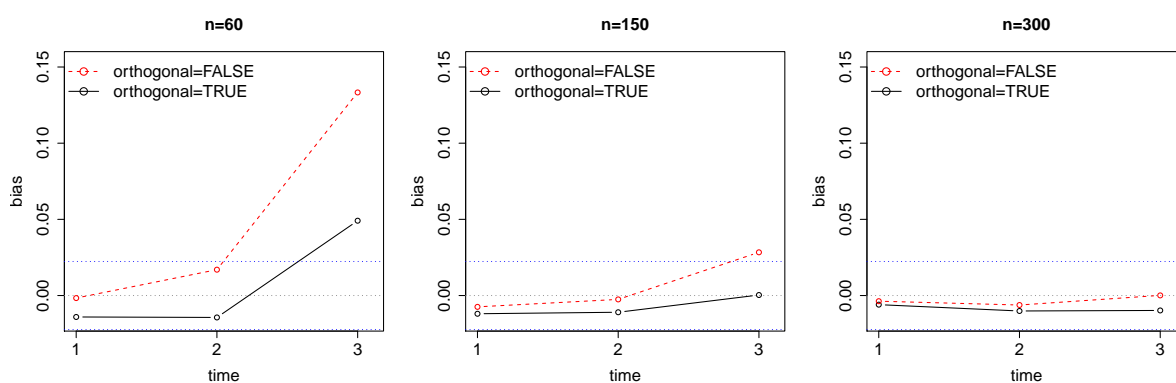


Figure E3. Average biases of the two-stage G-estimators, either assuming an orthogonal covariance matrix for the latent variables ('orthogonal=TRUE'), or allowing for correlated latent variables ('orthogonal=FALSE'), in the fitted measurement model. The true factor loadings depended on the occasion t . Both the mediator and outcome were continuous variables that were latent at all occasions. The true values of the direct effects were 0.64 ($t = 1$), 1.34 ($t = 2$), 2.80 ($t = 3$). The Monte Carlo sampling errors are drawn as horizontal blue dotted lines. The sample sizes n are stated in each panel.

Appendix F

Linear model for variables in Figure 2

The conditional mean models for the variables depicted in the right diagram of Figure 2 at times $t = 1, \dots, 3$ were:

$$E[L_t|X, \bar{L}_{t-1}, \bar{M}_{t-1}, \bar{Y}_{t-1}, U] = i_{lt} + q_t X + \sum_{t'=1}^{t-1} (f_{t't} L_{t'} + k_{t't} M_{t'} + s_{t't} Y_{t'}) + f_{0t} U \quad (5)$$

$$E[M_t|X, \bar{L}_t, \bar{M}_{t-1}, \bar{Y}_{t-1}, U_m] = i_{mt} + a_t X + g_{tt} L_t + \sum_{t'=1}^{t-1} (g_{t't} L_{t'} + d_{t't} M_{t'} + r_{t't} Y_{t'}) + d_{0t} U_m \quad (6)$$

$$E[Y_t|X, \bar{L}_t, \bar{M}_t, \bar{Y}_{t-1}, U_y, U] = i_{yt} + c_t X + b_{tt} M_t + h_{tt} L_t + \sum_{t'=1}^{t-1} (b_{t't} M_{t'} + h_{t't} L_{t'} + e_{t't} Y_{t'}) + e_{0t} U_y + h_{0t} U. \quad (7)$$

Appendix G

Bootstrap parameter estimates in analysis of randomised pilot study

In the analysis of the randomised eHealth pilot study, when motivation following treatment was assumed to be latent, about 10% of the bootstrap samples (used to estimate the standard errors of the controlled direct effects) had scores with negative measurement error variance estimates in the first stage. For such samples, the measurement error variances were set to zero. Histograms of the resulting measurement error variance estimates are plotted in Figure G1.

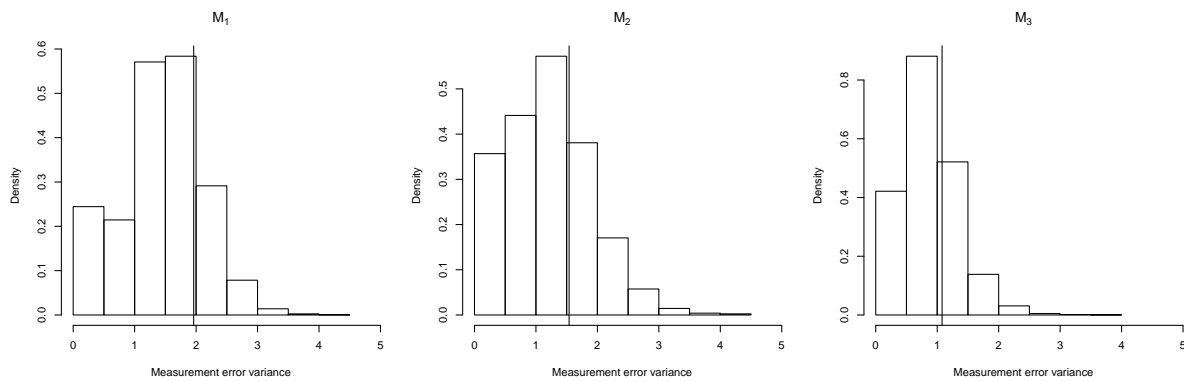


Figure G1. Histograms of bootstrap estimates of the measurement error variances of the mediator scores \widehat{M}_t at each time $t = 1, 2, 3$ using data from the randomised eHealth pilot study. The observed values are plotted as vertical solid lines.

References

- Browne, M. W. (1984). Asymptotically distribution-free methods for the analysis of covariance structures. *British Journal of Mathematical and Statistical Psychology*, *37*(1), 62–83. doi: 10.1111/j.2044-8317.1984.tb00789.x
- Efron, B., & Tibshirani, R. J. (1994). *An introduction to the bootstrap*. Chapman and Hall/CRC. doi: 10.1007/978-1-4899-4541-9
- Rosseel, Y. (2012). lavaan: An R package for structural equation modeling. *Journal of Statistical Software*, *48*(2), 1–36. doi: 10.18637/jss.v048.i02
- Vansteelandt, S., & Sjölander, A. (2016). Revisiting G-estimation of the effect of a time-varying exposure subject to time-varying confounding. *Epidemiologic Methods*, *5*(1), 37–56. doi: 10.1515/em-2015-0005