

# Causal Mediation Analysis

Tyler J. VanderWeele, Ph.D.

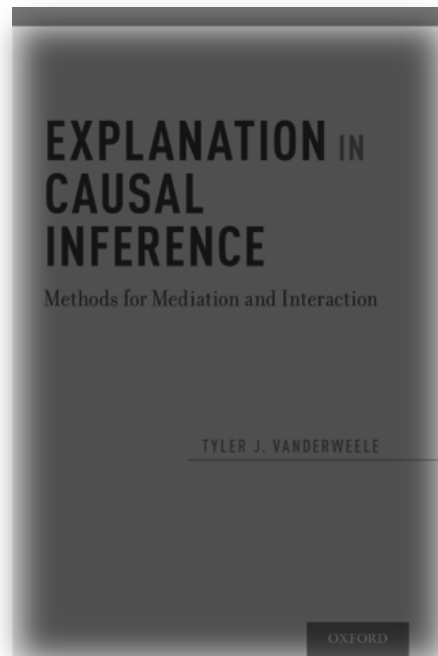
*Upcoming Seminar:*  
October 23-24, 2020, Remote Seminar

OXFORD UNIVERSITY PRESS

**Explanation in Causal  
Inference**  
Methods for Mediation and  
Interaction

2015 | Hardcover | ISBN: 9780199325870

Available on Amazon or from  
Oxford University Press



## Plan of Presentation

- (1) Concepts and Methods for Mediation
- (2) Sensitivity Analysis
- (3) Mediation with Time-to-Event Outcome
- (4) Multiple Mediators
- (5) Surrogate Outcomes
- (6) A Unification of Mediation and Interaction

# 1. Concepts and Methods

Tyler J. VanderWeele  
Departments of Epidemiology and Biostatistics  
Harvard T.H. Chan School of Public Health

## Plan of Presentation

- (1) Motivating Examples
- (2) Traditional Approaches and Limitations
- (3) Counterfactual Concepts
- (4) Regression-Based Methods
- (5) Binary Outcomes and Mediators
- (6) Empirical Examples
- (7) Macros and Software
- (8) Monte-Carlo Approach
- (9) Study Design

## Questions of Mediation

In a number of research contexts we might be interested in the extent to which the effect of some exposure A on some outcome Y is mediated by an intermediate variable M and to what extent it is direct



Stated another way, we are interested in the direct and indirect effects of the exposure

## Genetics Example

**Lung Cancer**: In 2008, three GWAS studies (Thorgeirsson et al., 2008; Hung et al., 2008; Amos et al., 2008) identified variants on chromosome 15q25.1 that were associated with lung cancer

**Smoking**: These variants had also been shown to be associated with smoking behavior (average cigarettes per day) e.g. through nicotine dependence (Saccone et al., 2007; Spitz et al., 2008)

**Debate**: there was debate as to whether the effect on lung is direct or operates through pathways related to smoking behavior (Chanock and Hunter, 2008); two thought direct, one mediated

**Interaction**: Complicating matter further there was some evidence gene-environment interaction: carriers of the variant allele extract more nicotine and toxins from each cigarette (Le Marchand, 2008)

## Perinatal Epidemiology Example

**ART**: There is evidence that use of assisted reproductive technologies (ART) lead to worse birth outcomes

**Twins**: It is also clear that use of ART leads to high incidence of twins; being born as a twin also leads to worse birth outcomes

**Mediation**? To what extent is the effect of ART on birth outcomes due to twinning? To what extent is it through other pathways?

**Policy Relevance**: Twins could mostly be prevented for those using ART by e.g. only allowing single embryo transfer

Some countries have adopted this policy e.g. Sweden

## Standard Approach

The standard approach to mediation analysis in much epidemiologic and social science research consists first of regressing the outcome Y on the exposure A and confounding factors C

$$E[Y|A=a,C=c] = \phi + \phi_1 a + \phi_2 c$$

And compare the estimate  $\phi$  of exposure A with the estimate  $\theta$  obtained when including the potential mediator M in the regression model

$$E[Y|A=a,M=m,C=c] = \theta + \theta_1 a + \theta_2 m + \theta_3 c$$

If the coefficients  $\phi$  and  $\theta$  differ The usual measures of direct and indirect effect then some of the effect is thought to be mediated and the following estimates are often used:

$$\text{Indirect effect} = \phi - \theta$$

$$\text{Direct effect} = \theta$$

## Standard Approach

Using the difference between the two coefficients is sometimes called the “difference method”

Another standard method, used more commonly in the social sciences is sometimes referred to as the “product method” (Baron and Kenny, 1986):

One regresses M on A:  $E[M|A=a,C=c] = \beta_0 + \beta_1 a + \beta_2 c$

One regresses Y on M and A:  $E[Y|A=a,M=m,C=c] = \theta_0 + \theta_1 a + \theta_2 m + \theta_3 c$

The direct effect is once again  $\theta_1$

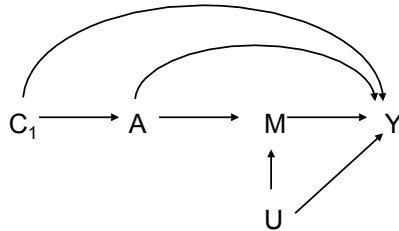
The indirect or mediated effect is the product of the coefficient of A in the regression for M times the coefficient of M in the regression for Y:  $\beta_1 \theta_2$

The product method and difference method will coincide for continuous outcomes provided the models are correctly specified but not for binary outcomes (MacKinnon and Dwyer, 1993, MacKinnon et al., 1995)

## Standard Approach

The standard approach to mediation analysis of just including the mediator in the regression is subject to two important limitations

PROBLEM 1: Even if the exposure is randomized or if all of the exposure-outcome confounders are included in the model there may be confounders of the mediator-outcome relationship



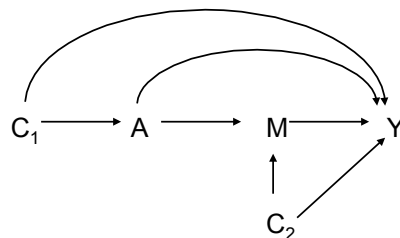
If control is not made for the mediator-outcome confounders then results from the standard approach can be highly biased

## Mediator-Outcome Confounding

In many social and biomedical studies, careful thought is given to control for confounding of the exposure-outcome relationship; data are collected on all variables thought to confound the relationship between the exposure and the outcome ( $C_1$  in the diagram)

However, often little thought is given to collecting data on variables that might confound the mediator-outcome relationship ( $C_2$  in the diagram)

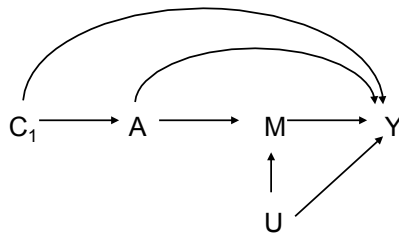
Mediation analyses are often secondary analyses in social and biomedical research and these variables often are not controlled for



## Mediator-Outcome Confounding

Just as unmeasured exposure-outcome confounders can generate confounding bias of estimates of overall effects

So also unmeasured mediator-outcome confounders can generate bias of estimates of direct and indirect effects



## Mediator-Outcome Confounding

The importance of controlling for mediator-outcome confounders when examining direct and indirect effects was also pointed out early on in the psychology literature on mediation (Judd and Kenny, 1981)

However a later paper in the psychology literature (Baron and Kenny, 1986) came to be the canonical reference for mediation analysis in the social sciences (>45,000 citations on Google Scholar)

Unfortunately, the Baron and Kenny (1986) paper did not note that control needed to be made for mediator-outcome confounders in the estimation of direct and indirect effects, even though the point had been made by Judd and Kenny five years earlier in 1981 and even though the two papers shared an author

As a result the point has been ignored by most of the research on mediation in the social sciences; many of these analyses are thus likely biased (possibly severely)

Contrary to claims made in the psychology literature, mediator-outcome <sup>14</sup> confounding is an issue for mediation analysis even in randomized trials!



## Mediator-Outcome Confounding

SMaRT trial (Strong et al., 2008): a randomized cognitive behavioral therapy intervention

Effect on depression symptoms after 3 months (SCL-20 depression, scale 0-4), was:

$$E[Y|A=1]-E[Y|A=0]=-0.34 \text{ (95\% CI: -0.55, -0.13)}$$

Intervention also had an effect on the use of antidepressant, M, at three months:

$$E[M|A=1]-E[M|A=0]=0.27$$

Those in the CBT arm were more likely to use antidepressants

Does the CBT intervention affect depressive symptoms simply because of higher antidepressant use, or other pathways?

What happens when we regress outcome Y on treatment and mediator (anti-depressant use)...?

## Mediator-Outcome Confounding

The coefficient for antidepressant use is positive!

It looks like antidepressant use increases depression!

The mediated effect through antidepressant use looks detrimental!

The “direct effect” looks larger than the total effect!

What is going on here...? Mediator-Outcome Confounding

Those in more difficult situations both use an antidepressant and have higher levels of depressive symptoms

When we ignore this confounding we get paradoxical results!

Anti-depressant use is not randomized here

But there are of course other trials that have randomized anti-depressant use...

## Mediator-Outcome Confounding

There are essentially two approaches to address mediator-outcome confounding (ideally both will be used):

(1) If mediation analysis is going to be part of an epidemiologic study then careful thought should be given to collecting data on mediator-outcome confounding variables during the study design stage

(2) After the study is finished, if there are unmeasured mediator-outcome confounders then sensitivity analysis techniques can be used to assess the extent to which the unmeasured confounding variable would have to affect the mediator and the outcome (and possibly the exposure) in order to invalidate inferences about direct and indirect effects (VanderWeele, 2010; Imai et al. 2010; Hafeman, 2011; Tchetgen Tchetgen and Shpitser, 2012)

## Exposure-Mediator Interactions

Limitation 2: Interactions between the effects of the exposure and the mediator, if present, and neglected can lead to biases

Even if we include an interaction term in the regression model:

$$E[Y|A=a,C=c] = \phi_0 + \phi_1 a + \phi_2' c$$

$$E[Y|A=a,M=m,C=c] = \theta_0 + \theta_1 a + \theta_2 m + \theta_3 a m + \theta_4' c$$

The usual measures of direct and indirect effect

$$\text{Indirect effect} = \phi_1 - \theta_1$$

$$\text{Direct effect} = \theta_1$$

break down because it is unclear how to handle the interaction coefficient  $\theta_3$

## Exposure-Mediator Interactions

In addition to clarifying the various no-unmeasured confounding assumptions that are needed in mediation analysis, the early causal inference literature on mediation (Robins and Greenland, 1992; Pearl, 2001) provided definitions of direct and indirect effects that could be used even when there were interaction between the effects of the exposure and the mediator on the outcome and that could also be used in the presence of non-linear models

In what follows we will:

- (1) Consider the causal (“counterfactual”) definitions of direct and indirect effects for mediation analysis and discuss the no-unmeasured confounding assumptions required for identification
- (2) Describe regression methods that can be used to estimate these counterfactual direct and indirect effect quantities (e.g. VanderWeele and Vansteelandt, 2009, 2010; cf. Imai et al., 2010)
- (3) Provide sensitivity analysis techniques to assess the importance of possible violations to the no-unmeasured confounding assumptions

## Definitions

Let  $Y$  denote some outcome of interest for each individual

Let  $A$  denote some exposure or treatment of interest for each individual

Let  $M$  denote some post-treatment intermediate(s) for each individual (potentially on the pathway between  $A$  and  $Y$ )

Let  $C$  denote a set of covariates for each individual

Let  $Y_a$  be the counterfactual outcome (or potential outcome)  $Y$  for each individual when intervening to set  $A$  to  $a$

Let  $M_a$  be the counterfactual outcome  $M$  for each individual when intervening to set  $A$  to  $a$

Let  $Y_{am}$  be the counterfactual outcome  $Y$  for each individual when intervening to set  $A$  to  $a$  and  $M$  to  $m$

## Definitions

Robins and Greenland (1992) and Pearl (2001) proposed the following counterfactual definitions for direct and indirect effects:

Controlled direct effect: The controlled direct effect comparing treatment level A=1 to A=0 intervening to fix M=m

$$\text{CDE}(m) = Y_{1m} - Y_{0m}$$

Natural direct effect: The natural direct effect comparing treatment level A=1 to A=0 intervening to fix M=M<sub>0</sub>

$$\text{NDE} = Y_{1M_0} - Y_{0M_0}$$

Natural indirect effect: The natural indirect effect comparing the effects of M=M<sub>1</sub> versus M=M<sub>0</sub> intervening to fix A=1

$$\text{NIE} = Y_{1M_1} - Y_{1M_0}$$

## Properties of Direct and Indirect Effects

A total effect decomposes into a direct and indirect effect:

$$\begin{aligned} Y_1 - Y_0 &= Y_{1M_1} - Y_{0M_0} \\ &= (Y_{1M_1} - Y_{1M_0}) + (Y_{1M_0} - Y_{0M_0}) \\ &= \text{NIE} + \text{NDE} \end{aligned}$$

The definitions of natural direct and indirect effect do not presuppose no interactions between the effects of the exposure and the mediator on the outcome

The effect decomposition of a total effect into a natural direct and indirect effect also does not presuppose no interaction between the effects of the exposure and the mediator on the outcome

Natural direct and indirect effects are useful for effect decomposition; in general, controlled direct effects are not

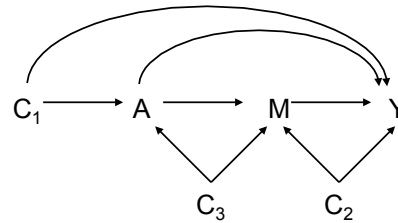
## Identification of Direct and Indirect Effects

To estimate average natural direct and indirect effects we need:

- (1) There are no unmeasured exposure-outcome confounders given C
- (2) There are no unmeasured mediator-outcome confounders given (C,A)
- (3) There are no unmeasured exposure-mediator confounders given C
- (4) There is no mediator-outcome confounder affected by exposure (i.e. no arrow from A to  $C_2$ )

For controlled direct effects, only assumptions (1) and (2) are needed

Note (1) and (3) are guaranteed when treatment is randomized



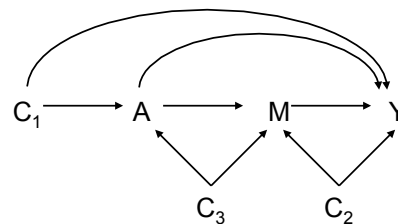
## Identification of Direct and Indirect Effects

More formally, in counterfactual notation, these assumptions are:

- (1) is  $Y_{am} \perp\!\!\!\perp A \mid C$
- (2) is  $Y_{am} \perp\!\!\!\perp M \mid C, A$
- (3) is  $M_a \perp\!\!\!\perp A \mid C$
- (4) is  $Y_{am} \perp\!\!\!\perp M_{a^*} \mid C$

For controlled direct effects, only assumptions (1) and (2) are needed

Note (1) and (3) are guaranteed when treatment is randomized



## Identification of Direct and Indirect Effects

Under assumptions (1) and (2) the controlled direct effect conditional on the covariates is given by:

$$E[\text{CDE}(m) | c] = E[Y|A=1, m, c] - E[Y|A=0, m, c]$$

Under (1)-(4) the conditional natural direct and indirect effects are:

$$E[\text{NDE} | c] = \sum_m \{E[Y|A=1, m, c] - E[Y|A=0, m, c]\} P(M=m|A=0, c)$$

$$E[\text{NIE} | c] = \sum_m E[Y|A=1, m, c] \{P(M=m|A=1, c) - P(M=m|A=0, c)\}$$

These are the effects within strata of the covariates

We could take averages over each stratum weighted by the probability  $P(C=c)$  to get population averages of the effects

## Regression for Causal Mediation Analysis

Similar concepts apply to treatment levels  $A=a$  to  $A=a^*$  (replace 1 by  $a$  and 0 by  $a^*$ )

Under our confounding assumptions (1)-(4), natural direct and indirect effects are given by the following expressions:

$$E[Y_{aM_{a^*}} - Y_{a^*M_{a^*}}] = \sum_{c,m} \{E[Y|a, m, c] - E[Y|a^*, m, c]\} P(m|a^*, c) P(c)$$

$$E[Y_{aM_a} - Y_{aM_{a^*}}] = \sum_{c,m} E[Y|a, m, c] \{P(m|a, c) - P(m|a^*, c)\} P(c)$$

We could consider fitting a parametric regression model for  $Y$  and a parametric regression model for  $M$  and computing this analytically (VanderWeele and Vansteelandt, 2009, 2010; Valeri and VanderWeele, 2013)

Alternatively Imai et al. (2010) propose to use a broad class of parametric or semiparametric models for  $Y$  and  $M$  and then to use simulations to calculate natural direct and indirect effects using the formulas above and the standard errors for these effects by bootstrapping

## Regression for Causal Mediation Analysis

We use regressions that accommodate exposure-mediator interaction:

$$E[Y|A=a, M=m, C=c] = \theta_0 + \theta_1 a + \theta_2 m + \theta_3 a m + \theta_4' c$$

$$E[M|A=a, C=c] = \beta_0 + \beta_1 a + \beta_2' c$$

Under assumptions (1)-(4), and provided our models are correctly specified, we can combine the estimates from the two models to get the following formulas for direct and indirect effects, comparing exposure levels  $a$  and  $a^*$  (VanderWeele and Vansteelandt, 2009):

$$CDE(a, a^*; m) = (\theta_1 + \theta_3 m)(a - a^*)$$

$$NDE(a, a^*; a^*) = (\theta_1 + \theta_3(\beta_0 + \beta_1 a^* + \beta_2' E[C]))(a - a^*)$$

$$NIE(a, a^*; a) = (\theta_2 \beta_1 + \theta_3 \beta_1 a)(a - a^*)$$

If the conditional NDE were of interest then we would have:

$$E[Y_{aM^*} - Y_{a^*M^*} | C=c] = (\theta_1 + \theta_3(\beta_0 + \beta_1 a^* + \beta_2' c))(a - a^*)$$

## Regression for Causal Mediation Analysis

Note that if there is no interaction between the effects of the exposure and the mediator on the outcome so that  $\theta_3=0$  then these expressions reduce to:

$$CDE(a, a^*; m) = NDE(a, a^*; a^*) = \theta_1(a - a^*)$$

$$NIE(a, a^*; a) = \theta_2 \beta_1(a - a^*)$$

which are the expressions often used for direct and indirect effects in the social science literature (Baron and Kenny, 1986) – the “product method”

However, unlike the Baron and Kenny (1986) approach, this approach to direct and indirect effects using counterfactual definitions and estimates can be employed even in settings in which an interaction is present

Standard errors can be obtained using the delta method

Proportion mediated is just the indirect effect divided by the total effect

**SAS, Stata, and SPSS macros** (Valeri and VanderWeele, 2013) can do this automatically for continuous, binary, count, and time-to-event outcomes