

# Analytical Uncertainty Quantification for Multilevel Mediation Analysis

July 11, 2024

## 1 Notes to Self

- In the paper, I organize variables as  $Y, X, M, W$ . This is based on defining counterfactuals as  $Y_x$ , then adding  $M$ , giving  $Y_{x,m}$ . In my code, it's always  $Y, M, X, W$ . I must be very careful to not mess this ordering up. The code is more arbitrary, but it's also harder to change. For now, I'm just going to try to juggle this difference manually, but I will need to carefully validate any results when passing from code to paper (this sort of issue has already been the source of one elusive error in my code).

## 2 Introduction

- Literature review
  - Samoilenko and Lefebvre
  - Imai et al.
  - SEM world
- Overview of our contribution
  - Analytical UQ for mediation effects
  - Existing work is Monte Carlo-based. E.g. Imai's method, bootstrap

Mediation analysis is a central problem in modern causal inference. Many scientific and public health questions have the form of separating the direct effect of some exposure on an outcome from the indirect effect of that exposure via a mediator. Many authors have developed methods to address this problem; gradually increasing in complexity. Early work by Baron and Kenny [1986] laid the foundation for future developments, but did not use the machinery of counterfactual outcomes.

There have been several approaches to the analysis of causal mediation. One group established a non-parametric identification result [Imai et al., 2010b], and

used this result to estimate mediation effects in various contexts [Imai et al., 2010a, 2011]. Their methodology is implemented in the R package `mediation` [Tingley et al., 2014]. Note that this group only estimates mediation effects as expected differences in counterfactuals, so if the outcome is binary then mediation effects are only available on the risk-difference scale.

Another approach began with the work of VanderWeele and Vansteelandt [2009] on continuous outcomes, and was later extended to handle various modifications to the basic model [VanderWeele and Vansteelandt, 2010, 2013, VanderWeele, 2014]. Of particular interest to us is the modification to handle binary outcomes [VanderWeele and Vansteelandt, 2010]. Effects are defined on the odds-ratio scale, and the outcome is assumed to be rare. Later work by Samoilenko and Lefebvre [2021] removes the rare-outcome assumption and extends the work of VanderWeele and Vansteelandt [2009] to handle effects on the risk-difference, risk-ratio and odds-ratio scales. See also Samoilenko et al. [2018], Samoilenko and Lefebvre [2023] for more details.

Briefly, causal mediation analysis is based on the counterfactual, or potential outcome framework. Let  $Y$  be an outcome of interest and  $X$  be an exposure. We write  $Y(x)$  for the value  $Y$  would have attained if, possibly counter to fact,  $X$  had been set to the value  $x$ . Introducing a mediator,  $M$ , we write  $M(x)$  for the value  $M$  would have attained if  $X = x$ , and  $Y(x, m)$  for the value of  $Y$  when  $X = x$  and  $M = m$ . Note that every individual in the population has values for each of the above quantities,  $Y(x)$ ,  $M(x)$ ,  $Y(x, m)$  at every possible value of  $x$  and  $m$ . Unfortunately, in practice we only observe  $Y$  and  $M$  for the values of  $x$  and  $m$  which actually occurred. This is known as the “fundamental problem of causal inference” [Ding and Li, 2018, Holland, 1986].

The standard approach to solving this fundamental problem is to avoid estimating individual-level counterfactuals and instead estimate population averages thereof. Under standard assumptions, such as consistency and no unmeasured confounders [see, e.g., Pearl, 2009], we can estimate expected counterfactuals as functions of conditional expectations. From this point, mediation analysis reduces to a problem of classical statistics; one which can be solved using traditional regression methodology. The difference between continuous and binary outcomes (or mediators) is essentially addressed by choosing between linear and logistic regression. Other data types (e.g., count, survival), or more flexible relationships (e.g., splines), can also be incorporated by selecting the appropriate regression methodology.

One extension which is of particular interest is to dependent data via multilevel, or mixed-effects, models. Mixed-effects regression models involve the introduction of random, unobserved coefficients to an existing regression [awk]. A common setting in which such a model arises is clustered data, where the random coefficients differ across clusters, but are constant within each cluster. See, e.g., Demidenko [2004] for an overview of mixed-effects methodology. When applied to mediation analysis, mixed-effects methods allow for modelling group-specific mediation effects, and the effect of this heterogeneity on the estimation of global effects.

In this paper, we present a general framework for multilevel mediation analy-

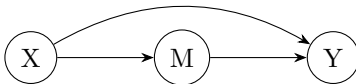


Figure 1: Causal diagram showing  $M$  mediating the effect of  $X$  on  $Y$ .

sis based on the estimation of nested counterfactuals. In particular, our method can be applied to estimate mediation effects on whatever scale is of interest (e.g., risk difference, risk ratio, odds ratio).

### 3 Multilevel Mediation Analysis

- Define counterfactuals (ctfs) and nested counterfactuals
- Define mediation effects
  - Continuous outcome
  - Binary outcome on risk difference, risk ratio and odds ratio scales
- Identification of mediation effects
  - See Imai et al. [2010a] for non-parametric identification of nested counterfactuals
- Regression modelling
  - Restrict to binary outcome?
  - Fixed effects regression (brief)
  - Mixed effects regression
  - Prediction of random effects

#### 3.1 Counterfactual-Based Mediation Effects

Our approach to causal mediation analysis is based on the counterfactual framework of ????. Briefly, let  $Y$  be an outcome of interest,  $X$  be an exposure which is a causal driver of  $Y$ , and  $M$  be a mediator, which influences  $Y$  and is influenced by  $X$ . Figure 1 shows a causal diagram representing this relationship. We call the top arrow the direct effect of  $X$  on  $Y$ , and the bottom path through  $M$  the indirect effect of  $Y$ . Taken together, these two pathways constitute the total effect of  $X$  on  $Y$ .

We define the counterfactual  $Y_x$  as the value  $Y$  would assume when  $X = x$ . Similarly, write  $M_x$  for the value  $M$  would assume when  $X = x$ . Next, write  $Y_{x,m}$  for the value  $Y$  would assume if  $X$  and  $M$  were set to  $x$  and  $m$  respectively. Combining these ideas, we get the “nested counterfactual”  $Y_{x,M_{x'}}$ , which is the value of  $Y$  when  $X = x$  and  $M$  is set to whatever it would have been when  $X = x'$ . It is common to link ordinary and nested counterfactuals by making

the consistency assumption, which states  $Y_x = Y_{x,M_x}$ . Note that if  $x \neq x'$  then the nested counterfactual is necessarily unobservable. It is nevertheless possible, given certain assumptions, to estimate the expected values of ordinary and nested counterfactuals. Section 3.2 goes into detail on the identification of expected counterfactuals with estimable quantities.

There are three main types of mediation effect: total, direct and indirect, although the scale on which these are measured can vary. For concreteness, we focus here only on discrepancies defined as differences; i.e., those of the form  $Y_{x_1,M_{x'_1}} - Y_{x_2,M_{x'_2}}$ . Extending our results to effects defined on different scales (e.g. ratios) is straightforward. We will return briefly to this in Section ???.

Proceeding now to the actual definitions, we define the total effect of  $X$  on  $Y$  to be  $TE(x, x') = EY_x - EY_{x'}$ . The direct effect is defined as  $DE(x, x') = Y_{x,M_{x'}} - Y_{x',M_{x'}}$ , and the indirect effect is  $IE(x, x') = Y_{x,M_x} - Y_{x,M_{x'}}$ . See, e.g., ??? for motivation and discussion of these definitions.

### 3.2 Identification and Modelling of Expected Counterfactuals

A fundamental problem of causal inference is that we can only ever observe one counterfactual outcome on a particular individual. In mediation analysis, this problem is even worse, since many of our definitions involve the nested counterfactual,  $Y_{x,M_{x'}}$ , which when  $x \neq x'$  cannot be observed on any individual. Nevertheless, Imai et al. [2010a] give conditions under which the population average of a nested counterfactual can be expressed in terms of conditional expectations, possibly conditional on one or more additional covariates,  $W$ . Specifically, their Theorem 1 states that, under a condition they call “Sequential Ignorability”, we can write

$$E(Y_{x,M_{x'}}|W = w) = E_M[E_Y(Y|X = x, M, W = w)|X = x', W = w]. \quad (1)$$

In fact, their Theorem 1 is somewhat more general, giving an expression for the density of the nested counterfactual rather than its expected value.

Using Equation (1), we can estimate expected nested counterfactuals by working with the more tractable conditional expectations of  $Y$  and  $M$ . We model the latter using regression, either linear or logistic depending on the forms of  $Y$  and  $M$ . For concreteness, we take  $Y$  and  $M$  to both be binary. The extension of our method to continuous outcome and/or mediator is straightforward (comment on sums  $\rightarrow$  integrals and quadrature?). In this case, Equation (1) has a particularly simple form:

$$E(Y_{x,M_{x'}}|W = w) = \mathbb{P}(Y = 1|X = x, M = 1, W = w) \mathbb{P}(M = 1|X = x', W = w) + \mathbb{P}(Y = 1|X = x, M = 0, W = w) \mathbb{P}(M = 0|X = x', W = w). \quad (2)$$

We estimate the conditional probabilities on the right-hand side of Equation (2) using logistic regression. Note that we must fit two models: one to predict  $M$  using  $X$  and  $W$ , and another to predict  $Y$  using  $M$ ,  $X$  and  $W$ .

We extend this regression modelling with the introduction of random effects [see, e.g. Demidenko, 2004]. Specifically, we include random effects for the intercept and  $X$  in our model for  $M$ , and for the intercept,  $X$  and  $M$  in our model for  $Y$ . Let  $U$  and  $V$  be the random effects for our models of  $Y$  and  $M$  respectively. We can re-write Equation (2) as

$$\begin{aligned}\mathbb{E}(Y_{x,M_{x'}}|W=w) &= [\mathbb{E}_U \mathbb{P}(Y=1|U, X=x, M=1, W=w) \cdot \\ &\quad \mathbb{E}_V \mathbb{P}(M=1|V, X=x', W=w)] + \\ &\quad [\mathbb{E}_U \mathbb{P}(Y=1|U, X=x, M=0, W=w) \cdot \\ &\quad \mathbb{E}_V \mathbb{P}(M=0|V, X=x', W=w)].\end{aligned}\quad (3)$$

Following the usual approach, we model the random effects  $U$  and  $V$  as normally distributed, allowing for correlation between effects from the same model, but assuming independence between models. Let  $U \sim N(0, \Gamma_Y)$  and  $V \sim N(0, \Gamma_M)$ .

### 3.2.1 Expanding One Term in Equation (3)

We derive an expression for the first term in Equation (3), then report results for the other three terms (awk).

First, write  $\eta_Y = (\alpha_0 + \alpha_X x + \alpha_M m + A_W^T w) + (U_0 + U_X x + U_M m)$  for the linear predictor of  $Y$  based on  $X$ ,  $M$  and  $W$ . Let  $\mu_Y = \alpha_0 + \alpha_X x + A_W^T w$ , so that the fixed-effects component of  $\eta_Y$  is  $\mu_Y + \alpha_M m$ , and let  $\xi_Y = U_0 + U_X x + U_M m$  be the random-effects component of  $\eta_Y$ . For convenience, we will also write  $\gamma_Y^2(c_1, c_2, c_3) = (c_1, c_2, c_3) \Gamma_Y(c_1, c_2, c_3)^T$ , so that  $\mathbb{V}\xi_Y = \gamma_Y^2(1, x, m)$  (recall that  $\Gamma_Y$  is the covariance matrix of  $U$ ).

It is a well-known fact about logistic regression that

$$\mathbb{P}(Y=1|U, X=x, M=1, W=w) = [1 + \exp(-\eta_Y)]^{-1}, \quad (4)$$

so the first term in Equation (3) can be written as

$$\mathbb{E}_U \mathbb{P}(Y=1|U, X=x, M=1, W=w) = \int \frac{\phi_3(u; 0, \Gamma_Y)}{1 + \exp(-\eta_Y)} du, \quad (5)$$

where  $\phi_d$  is the  $d$ -variate normal density. A straightforward change of variables gives us the alternative expression

$$\mathbb{E}_U \mathbb{P}(Y=1|U, X=x, M=1, W=w) = \int \frac{\phi_1(z; 0, 1)}{1 + \exp(-\mu_Y - \alpha_M - \gamma_Y(1, x, 1)z)} dz. \quad (6)$$

Importantly, the integral in Equation (6) is univariate (and thus amenable to numerical evaluation using quadrature). This integral arises often enough that we give it a name. Let

$$\Psi(a, b) = \int \frac{\phi(z; 0, 1)}{1 + \exp(-a - bz)} dz. \quad (7)$$

We can now write the first term in (3) compactly as

$$\mathbb{E}_U \mathbb{P}(Y = 1 | U, X = x, M = 1, W = w) = \Psi(\mu_Y + \alpha_M, \gamma_Y(1, x, 1)) \quad (8)$$

Returning now to the other terms in Equation (3), similar expressions hold. Write  $\eta_M = (\beta_0 + \beta_X x + B_W^T w) + (V_0 + V_X x)$  for the linear predictor of  $M$  based on  $X$  and  $W$ . Write  $\mu_M = \beta_0 + \beta_X x + B_W^T w$  and  $\xi_M = V_0 + V_X x$  for the fixed and random components respectively of  $\eta_M$ . Finally, write  $\gamma_M^2(c_1, c_2) = (c_1, c_2) \Gamma_M (c_1, c_2)^T$ , so that  $\mathbb{V} \xi_M = \gamma_M^2(1, x)$  ( $\Gamma_M$  is the covariance matrix of  $V$ ). We can now write the expected counterfactual in (3) as

$$\begin{aligned} \mathbb{E}(Y_{x, M_{x'}} | W = w) = & [\Psi(\mu_Y + \alpha_M, \gamma_Y(1, x, 1)) \cdot \Psi(\mu_M, \gamma_Y(1, x'))] + \\ & [\Psi(\mu_Y, \gamma_Y(1, x, 0)) \cdot \Psi(-\mu_M, \gamma_Y(1, x'))] \end{aligned} \quad (9)$$

Recall that each term in (9) is a univariate integral, and can thus be accurately evaluated using standard numerical quadrature routines available in most software packages.

### 3.3 Alternative Effect Definitions

This section needs to be completely re-written. The following might contain useable parts, but no guarantees.

When the outcome,  $Y$ , is continuous, it is natural to take the discrepancy function to be  $d(y_1, y_2) = y_1 - y_2$  as mentioned above. However, when  $Y$  is binary, it is less obvious which scale to use. In this setting, expected counterfactuals are probabilities, so mediation effects measure the discrepancy between two probabilities. One option is to use the same difference preferred for continuous outcome; we refer to this as the risk-difference scale. Alternatives include  $d(y_1, y_2) = y_1/y_2$ , the risk-ratio scale, and  $d(y_1, y_2) = [y_1/(1-y_1)]/[y_2/(1-y_2)]$ , the odds-ratio scale. We focus only on the binary outcome setting, although our results are easily extended to handle continuous outcomes. While our method can handle any choice of discrepancy function  $d$  (differentiable? Whatever is necessary to apply the  $\delta$ -method later; might need to be continuously differentiable), in the interest of accessibility we specialize our results to the three discrepancies just mentioned for binary outcomes.

### 3.4 Old

In this section, we define the total, direct and indirect mediation effects in terms of nested counterfactuals. We begin by applying the Mediation Formula of Pearl [2012] to identify nested counterfactuals with simple functions of conditional expectations.

We begin by observing that mediation effects are often defined in terms of nested counterfactuals (regardless of the scale on which these effects are reported). There is some causal inference theory to be done here. For now, I'm just going to write what I expect to be true, then later I will go back and add the necessary assumptions. As such, we begin by identifying a general

nested counterfactual with estimable quantities. To this end, write  $Y(x, m)$  for the counterfactual value of  $Y$  when  $X$  and  $M$  are set to  $x$  and  $m$  respectively. Similarly, write  $M(x)$  for the counterfactual value of  $M$  when  $X = x$ . We write  $Y(x, M(x'))$  for the value of  $Y$  when  $X$  is set to  $x$  and  $M$  is set to what it would have been if  $X$  were  $x'$ . We refer to  $Y(x, M(x'))$  as a nested counterfactual. Note that, when  $x \neq x'$ , the nested counterfactual is necessarily unobservable. Nevertheless, given regularity conditions, we can write

$$\mathbb{E}Y(x, M(x')) = \int \mathbb{E}(Y|M = m, X = x)\mathbb{P}(M = dm|X = x') \quad (10)$$

A similar expression holds conditional on pre-treatment covariates,

$$\mathbb{E}Y_c(x, M(x')) = \int \mathbb{E}(Y|M = m, X = x, C = c)\mathbb{P}(M = dm|X = x', C = c) \quad (11)$$

Having identified expected nested counterfactuals with integrals of conditional quantities, we can use regression to estimate the RHS of Equations 10 and 11. This identification with regression-based quantities is especially simple if  $M$  is binary. Here, we get  $\mathbb{P}(M = m|\cdot) = \mathbb{E}(M = m|\cdot)$ , where the latter quantity is popularly modelled with logistic regression. For simplicity, we hereafter assume that the mediator,  $M$ , is binary. **If not, we need to do something different. One option is to extract the conditional density of  $(M|\cdot)$ , then perform the integral. Another is to do regression with response  $\mathbb{E}(Y|M = m, X = x, C = c)$ , viewed as a function of  $m$ . There is likely an equivalence between these two methods; I can come back to this later.**

The above development allows only fixed-effects in the regression models. If we want to also incorporate random effects, we simply view Equations 10 and 11 as having had these random effects marginalized out. That is, if we write  $U$  and  $V$  for the random effects in our models for  $Y$  and  $M$  respectively, and  $G_U, G_V$  for the random effects' distributions, then we can re-write Equations 10 and 11 as

$$\mathbb{E}Y(x, M(x')) = \int \mathbb{E}(Y|U = u, M = m, X = x)\mathbb{P}(M = dm|V = v, X = x')G_U(du)G_V(dv) \quad (12)$$

$$\mathbb{E}Y_c(x, M(x')) = \int \mathbb{E}(Y|U = u, M = m, X = x, C = c)\mathbb{P}(M = dm|V = v, X = x', C = c)G_U(du)G_V(dv) \quad (13)$$

### 3.5 Ideas for CI Theory

Applying the Counterfactual Unnesting Theorem of Correa et al. [2021], and provided that the conditions specified therein are satisfied, we have

Alternatively, see Theorem 1 of Imai et al. [2010a]. This is specifically the result I want, not the corresponding identification results of Imai et al. [2010b], which focuses more narrowly on identification of mediation effects. Note that, in

Imai et al. [2010a], the function  $f$  is a density (with respect to some unspecified measure, possibly counting measure). In order to get the above formulas for expected values of  $Y$ , we do need to assume that  $\mathbb{E}Y$  is finite so we can apply Fubini’s Theorem (the one for finite integrals).

## 4 Estimation and Inference

- Estimate models using `lme4`
- UQ for model parameters (mention `merDeriv` package)
- UQ for nested counterfactuals via  $\delta$ -method
- Simultaneous UQ for triples of mediation effects (i.e. total, direct and indirect)
  - Can also do all 9 simultaneously
- UQ for predicted group-level effects?

## 5 Empirical Investigation

- Monte Carlo study
  - Proof of concept
  - Explore robustness. See Samoilenko and Lefebvre [2023] for inspiration.
- Real data
  - Trust study dataset?

## References

- Reuben M. Baron and David A. Kenny. The moderator-mediator variable distinction in social psychological research: conceptual, strategic, and statistical considerations. *Journal of Personality and Social Psychology*, 51(6), 1986.
- Juan D. Correa, Sanghack Lee, and Elias Bareinboim. Nested counterfactual identification from arbitrary surrogate experiments. *Advances in Neural Information Processing Systems*, 34, 2021.
- Eugene Demidenko. *Mixed models: theory and applications*. Wiley, 2004.
- Peng Ding and Fan Li. Causal inference: a missing data perspective. *Statistical Science*, 33(2), 2018.



- Paul W. Holland. Statistics and causal inference. *Journal of the American Statistical Association*, 81(396), 1986.
- Kosuke Imai, Luke Keele, and Dustin Tingley. A general approach to causal mediation analysis. *Psychological Methods*, 15(4), 2010a.
- Kosuke Imai, Luke Keele, and Teppei Yamamoto. Identification, inference and sensitivity analysis for causal mediation effects. *Statistical Science*, 25(1), 2010b.
- Kosuke Imai, Luke Keele, Dustin Tingley, and Teppei Yamamoto. Unpacking the black box of causality: learning about causal mechanisms from experimental and observational studies. *American Political Science Review*, 105(4), 2011.
- Judea Pearl. *Causality*. Cambridge University Press, 2nd edition, 2009.
- Judea Pearl. The causal mediation formula - a guide to the assessment of pathways and mechanisms. *Prevention Science*, 13(4), 2012.
- Mariia Samoilenko and Geneviève Lefebvre. Parametric-regression-based causal mediation analysis of binary outcomes and binary mediators: moving beyond the rareness or commonness of the outcome. *American Journal of Epidemiology*, 190(9), 2021.
- Mariia Samoilenko and Geneviève Lefebvre. An exact regression-based approach for the estimation of natural direct and indirect effects with a binary outcome and a continuous mediator. *Statistics in Medicine*, 42(3), 2023.
- Mariia Samoilenko, Lucie Blais, and Geneviève Lefebvre. Comparing logistic and log-binomial models for causal mediation analyses of binary mediators and rare binary outcomes: evidence to support cross-checking of mediation results in practice. *Observational Studies*, 4(1), 2018.
- Dustin Tingley, Teppei Yamamoto, Kentaro Hirose, Luke Keele, and Kosuke Imai. **mediation**: R package for causal mediation analysis. *Journal of Statistical Software*, 59(5), 2014.
- Tyler VanderWeele and Stijn Vansteelandt. Mediation analysis with multiple mediators. *Epidemiologic Methods*, 2(1), 2013.
- Tyler J. VanderWeele. A unification of mediation and interaction: a 4-way decomposition. *Epidemiology*, 25(5), 2014.
- Tyler J. VanderWeele and Stijn Vansteelandt. Conceptual issues concerning mediation, interventions and composition. *Statistics and its Interface*, 2(4), 2009.
- Tyler J. VanderWeele and Stijn Vansteelandt. Odds ratios for mediation analysis for a dichotomous outcome. *American Journal of Epidemiology*, 172(12), 2010.