# SOCI 40258

Causal Mediation Analysis
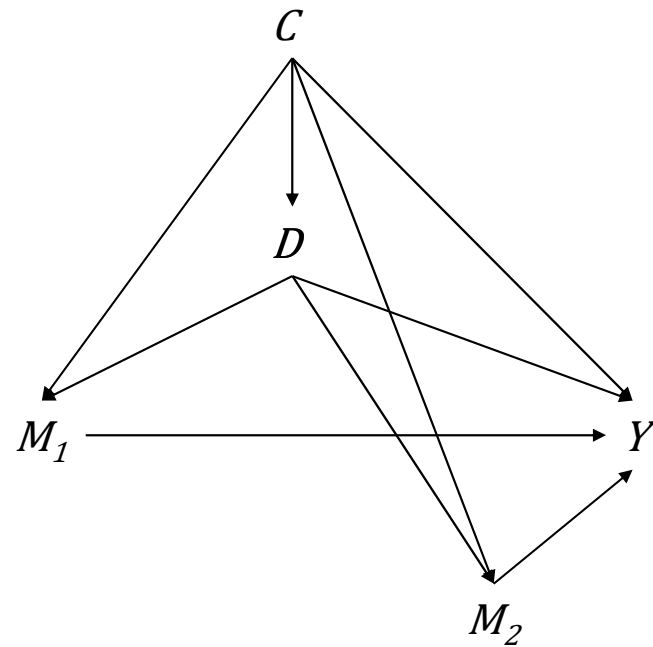
Week 7: Multiple Mediators

# Outline

- Graphical mediation models

- Natural effects through multiple mediators

- Nonparametric identification and estimation
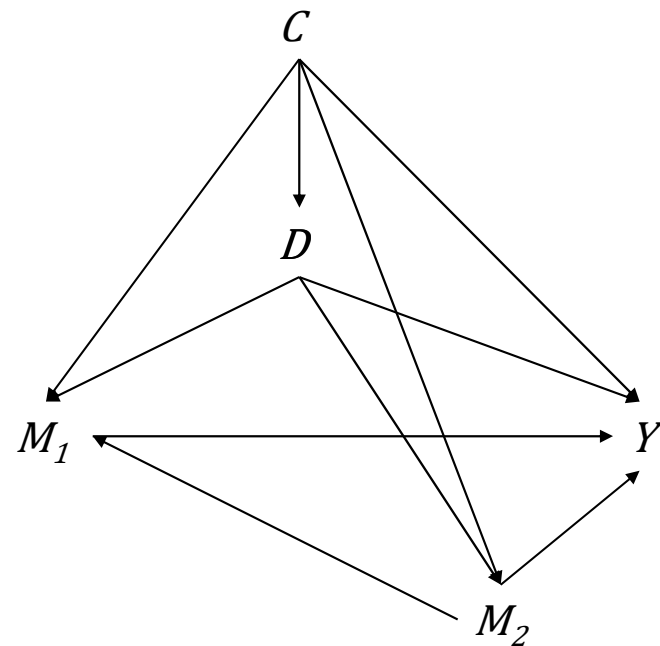
- Parametric estimation strategies

# Models with multiple mediators

- In this model, the exposure $D$ affects two mediators, $M_1$ and $M_2$, which both affect the outcome $Y$

- $M_1$ does not affect $M_2$, nor does $M_2$ affect $M_1$—that is, the two mediators are causally independent
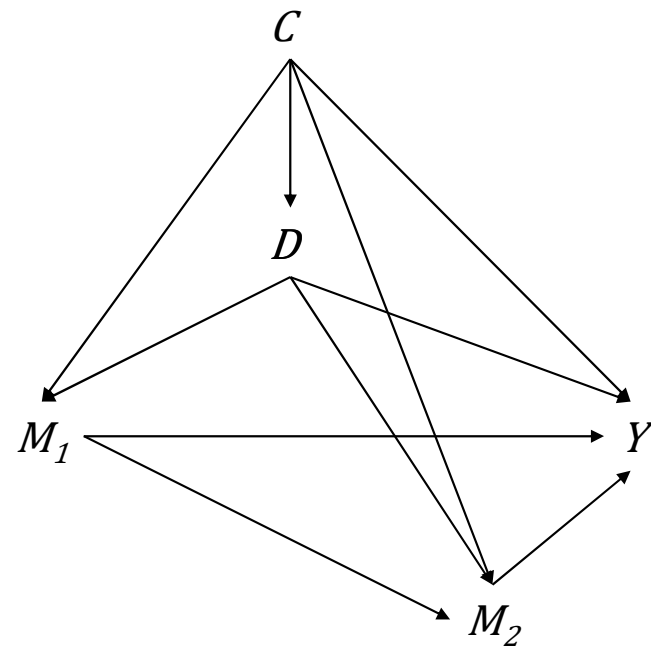
$C$

$D$

$M_1$

$Y$

$M_2$

# Models with multiple mediators

- In this model, the exposure $D$ affects two mediators, $M_1$ and $M_2$, which both affect the outcome $Y$

- $M_2$ now affects $M_1$, such that the mediators are causally dependent

- $M_2$ is an exposure-induced confounder with respect to the effect of $M_1$ on $Y$

# Models with multiple mediators

- In this model, the exposure $D$ affects two mediators, $M_1$ and $M_2$, which both affect the outcome $Y$

- $M_1$ now affects $M_2$, such that the mediators are again causally dependent

- $M_1$ is an exposure-induced confounder with respect to the effect of $M_2$ on $Y$

# Graphical mediation models

- The methods covered today are appropriate for data arising from a causal process resembling any of the graphical models depicted previously

- My presentation of these methods is tailored for models that allow general patterns of baseline confounding and causal dependence among the mediators

- These methods are also appropriate for settings without any baseline confounding and/or where the mediators are independent

# Natural effects with multiple mediators

- Natural effects with multiple mediators are very similar to the natural effects we have discussed previously, except they are defined in terms of a vector of $K$ mediators, denoted by $\mathbf{M} = \{M_1, M_2, \ldots, M_K\}$

- Specifically, with multiple mediators, the average total effect of the exposure on the outcome can be decomposed into direct and indirect components as follows:

$$ATE(d, d^*) = E\big(Y(d) - Y(d^*)\big)$$

$$= E\left(Y(d, \mathbf{M}(d)) - Y(d^*, \mathbf{M}(d^*))\right)$$

$$= E\left(Y(d, \mathbf{M}(d^*)) - Y(d^*, \mathbf{M}(d^*))\right) + E\left(Y(d, \mathbf{M}(d)) - Y(d, \mathbf{M}(d^*))\right)$$

# Natural effects with multiple mediators

- Natural effects with multiple mediators are very similar to the natural effects we have discussed previously, except they are defined in terms of a vector of $K$ mediators, denoted by $\mathbf{M} = \{M_1, M_2, \ldots, M_K\}$

- Specifically, with multiple mediators, the average total effect of the exposure on the outcome can be decomposed into direct and indirect components as follows:

$$ATE(d, d^*) = E\big(Y(d) - Y(d^*)\big)$$

$$= E\left(Y(d, \mathbf{M}(d)) - Y(d^*, \mathbf{M}(d^*))\right)$$

$$= \underbrace{E\left(Y(d, \mathbf{M}(d^*)) - Y(d^*, \mathbf{M}(d^*))\right)}_{\text{natural direct effect}} + \underbrace{E\left(Y(d, \mathbf{M}(d)) - Y(d, \mathbf{M}(d^*))\right)}_{\text{natural indirect effect}}$$

# The multivariate natural direct effect

- The multivariate natural direct effect:

$$MNDE(d, d^*) = E\left(Y(d, \mathbf{M}(d^*)) - Y(d^*, \mathbf{M}(d^*))\right)$$

$$= E\left(Y(d, M_1(d^*), \dots, M_K(d^*)) - Y(d^*, M_1(d^*), \dots, M_K(d^*))\right)$$

- The $MNDE(d, d^*)$ is the expected difference in the outcome if individuals had been exposed to $d$ rather than $d^*$ and if they had experienced the levels of all $K$ mediators that would have arisen naturally for them under exposure $d^*$

- It captures an effect of the exposure $D$ on the outcome $Y$ that operates through all mechanisms other than those involving the vector of mediators $\mathbf{M} = \{M_1, M_2, \dots, M_K\}$

# The multivariate natural direct effect

- The multivariate natural direct effect:

$$MNDE(d, d^*) = E\left(Y(d, \mathbf{M}(d^*)) - Y(d^*, \mathbf{M}(d^*))\right)$$

$$= E\left(Y(d, M_1(d^*), \ldots, M_K(d^*)) - Y(d^*, M_1(d^*), \ldots, M_K(d^*))\right)$$

- The $MNDE(d, d^*)$ isolates an effect not involving the mediators by…

  - comparing outcomes across different levels of the exposure ($d$ versus $d^*$)…

  - while holding all the mediators constant at their values under only one level of the exposure $\mathbf{M}(d^*) = \{M_1(d^*), M_2(d^*), \ldots, M_K(d^*)\}$

- This comparison deactivates the component of the total effect that is transmitted through all causal chains from exposure to the outcome operating through any of the mediators in $\mathbf{M}$

# The multivariate natural indirect effect

- The multivariate natural indirect effect:

$$MNIE(d, d^*) = E\left(Y(d, \mathbf{M}(d)) - Y(d, \mathbf{M}(d^*))\right)$$

$$= E\left(Y(d, M_1(d), \dots, M_K(d)) - Y(d, M_1(d^*), \dots, M_K(d^*))\right)$$

- The $MNIE(d, d^*)$ is the expected difference in the outcome if individuals had been exposed to $d$ and then…

  - experienced the levels of all mediators that would have arisen naturally for them under exposure $d$ rather than the levels that would have arisen naturally under exposure $d^*$

- It captures an effect of the exposure $D$ on the outcome $Y$ that operates through all mechanisms involving any of the mediators in $\mathbf{M}$

# The multivariate natural indirect effect

- The multivariate natural indirect effect:

$$MNIE(d, d^*) = E\left(Y\big(d, \mathbf{M}(d)\big) - Y\big(d, \mathbf{M}(d^*)\big)\right)$$

$$= E\left(Y\big(d, M_1(d), \ldots, M_K(d)\big) - Y\big(d, M_1(d^*), \ldots, M_K(d^*)\big)\right)$$

- The $MNIE(d, d^*)$ isolates an effect operating through all the mediators jointly by holding the exposure for each individual constant at $d\ldots$

  - while comparing outcomes across differences in all the mediators that would have arisen under different exposures, $\mathbf{M}(d)$ versus $\mathbf{M}(d^*)$

- This comparison deactivates all mechanisms from exposure to the outcome except for the causal chains operating through the vector of mediators

# Nonparametric identification

- Multivariate natural direct and indirect effects can be nonparametrically identified if the following conditions are met:

Assumption MNE.1: $Y(d, \mathbf{m}) \perp D | C$

Assumption MNE.2: $Y(d, \mathbf{m}) \perp \mathbf{M} | C, D = d$

Assumption MNE.3: $\mathbf{M}(d) \perp D | C$

Assumption MNE.4: $Y(d, \mathbf{m}) \perp \mathbf{M}(d^*) | C$

Assumption MNE.5: $P(d, \mathbf{m} | c) > 0$

Assumption MNE.6: $Y = Y(D) = Y\big(D, \mathbf{M}(D)\big) = Y(D, \mathbf{M})$

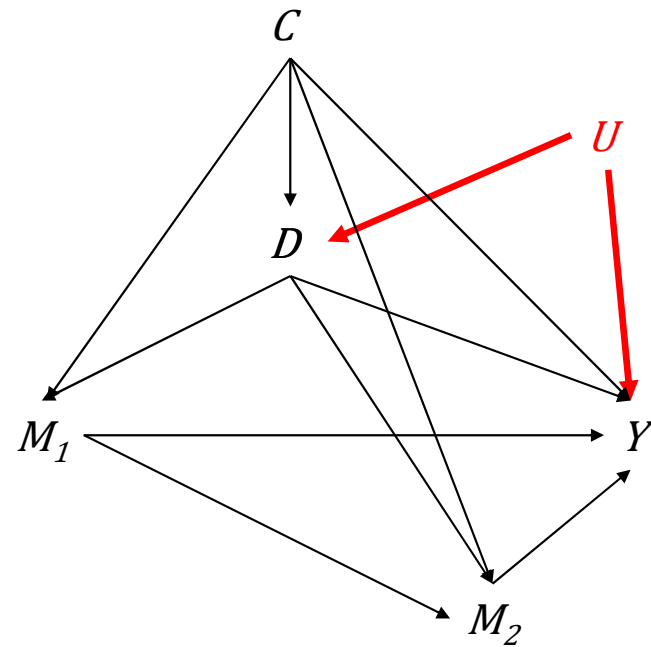# No unobserved exposure-outcome confounding

- Assumption MNE.1:

$$Y(d, \mathbf{m}) \perp D | C$$

where $Y(d, \mathbf{m}) = Y(d, m_1, \ldots, m_K)$

- This assumption requires that the exposure $D$ must be statistically independent of the joint potential outcomes $Y(d, \mathbf{m})$, conditional on the baseline confounders $C$

- Substantively, this assumption requires that there must not be any unobserved factors that confound the exposure-outcome relationship

# No unobserved exposure-outcome confounding

- Assumption MNE.1 would be violated if an unobserved variable jointly affects the exposure and outcome

- In this graph, $U$ is an unobserved confounder for the $D \rightarrow Y$ relationship

# No unobserved mediator-outcome confounding

- Assumption MNE.2:

$$Y(d, \mathbf{m}) \perp \mathbf{M} | C, D = d$$

- This assumption requires that the vector of mediators $\mathbf{M}$ must be statistically independent of the joint potential outcomes $Y(d, \mathbf{m})$, conditional on the baseline confounders $C$ in the group exposed to $d$

- Substantively, this assumption requires that there must not be any unobserved factors that confound the relationship between any one of the mediators and the outcome

# No unobserved mediator-outcome confounding

- Assumption MNE.2 would be violated if an unobserved variable jointly affects any one of the mediators and the outcome

- In this graph, $U_1$ is an unobserved confounder for the $M_1 \rightarrow Y$ relationship

- And $U_2$ is an unobserved confounder for the $M_2 \rightarrow Y$ relationship

# No unobserved exposure-mediator confounding

- Assumption MNE.3:

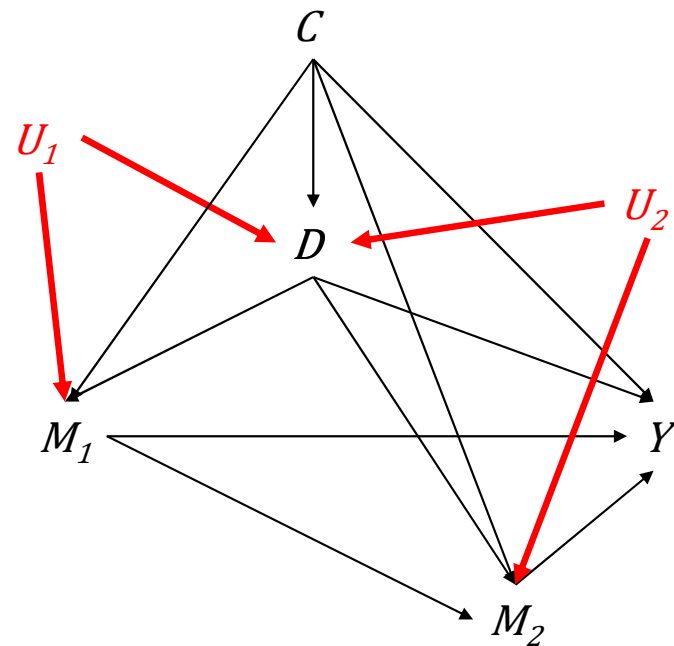  $$\mathbf{M}(d) \perp D | C$$

- This assumption requires that the exposure $D$ must be statistically independent of the potential values of all the mediators in $\boldsymbol{M}(d)$, conditional on the baseline confounders $C$

- Substantively, this assumption requires that there must not be any unobserved factors that confound the relationship between the exposure and any one of the mediators

# No unobserved exposure-mediator confounding

- Assumption MNE.3 would be violated if an unobserved variable jointly affects the exposure and any of the mediators

- In this graph, $U_1$ is an unobserved confounder for the $D \rightarrow M_1$ relationship

- And $U_2$ is an unobserved confounder for the $D \rightarrow M_2$ relationship

# No exposure-induced confounding

- Assumption MNE.4:

$$Y(d, \mathbf{m}) \perp \mathbf{M}(d^*) | C$$

- This assumption requires that the potential values of the mediators under exposure $d$ must be independent of the joint potential outcomes under exposure $d^*$, conditional on the baseline confounders $C$

- Known as a cross-world independence assumption, it requires that there must not be any exposure-induced confounders for any of the mediator-outcome relationships, whether they are observed or not

# No exposure-induced confounding

- Assumption MNE.4 would be violated if any variable, observed or not, jointly affects any of the mediators and the outcome…

  - …and is also affected by the exposure

- In this graph, $L$ is a confounder for the $M_2 \rightarrow Y$ relationship that is affected by $D$

# No exposure-induced confounding

- If $L$ is observed, it can be included in the vector of other mediators $\boldsymbol{M}$ and analyzed concurrently with them

- Including any exposure-induced confounders in $\boldsymbol{M}$ as additional mediators obviates violations of MNE.4

# Nonparametric identification

- Under assumptions MNE.1 to MNE.6, the multivariate natural direct effect can be equated with a function of observable data rather than nested and cross-world potential outcomes

- Nonparametric identification formula for the multivariate natural direct effect:

$$MNDE(d, d^*) = E\left(Y(d, \mathbf{M}(d^*)) - Y(d^*, \mathbf{M}(d^*))\right)$$

$$= \sum_c \sum_{\mathbf{m}} [E(Y|c, d, \mathbf{m}) - E(Y|c, d^*, \mathbf{m})] P(\mathbf{m}|c, d^*) P(c)$$

$$= \sum_c \sum_{\mathbf{m}} [E(Y|c, d, \mathbf{m}) - E(Y|c, d^*, \mathbf{m})] \prod_{k=1}^{K} P(m_k|c, d^*, \mathbf{m}_{k-1}) P(c)$$

where $\mathbf{m}_{k-1} = \{m_1, \dots, m_{k-1}\}$

# Nonparametric identification

- Under assumptions MNE.1 to MNE.6, the multivariate natural indirect effect can also be equated with a function of observable data rather than nested and cross-world potential outcomes

- Nonparametric identification formula for the multivariate natural indirect effect:

$$MNIE(d, d^*) = E\left(Y\big(d, \mathbf{M}(d)\big) - Y\big(d, \mathbf{M}(d^*)\big)\right)$$

$$= \sum_c \sum_{\mathbf{m}} E(Y|c, d, \mathbf{m})[P(\mathbf{m}|c, d) - P(\mathbf{m}|c, d^*)]P(c)$$

$$= \sum_c \sum_{\mathbf{m}} E(Y|c, d, \mathbf{m})\left[\prod_{k=1}^K P(m_k|c, d, \mathbf{m}_{k-1}) - \prod_{k=1}^K P(m_k|c, d^*, \mathbf{m}_{k-1})\right]P(c)$$

where $\mathbf{m}_{k-1} = \{m_1, \ldots, m_{k-1}\}$

# Nonparametric estimation

- Nonparametric identification involves equating causal effects defined in terms of counterfactuals with empirical quantities defined in terms of observable data, while ignoring random variability due to sampling

- In practice, however, we rarely have data from an entire target population and thus cannot simply ignore random variability due to sampling from this population

- Nonparametric estimation just involves plugging in sample analogs for the population quantities in the nonparametric identification formulas outlined previously

# Limitations of nonparametric estimation

- Limitations of nonparametric estimation

  - Sparsity

  - Curse of dimensionality

  - Excessive sampling variability

- With multiple mediators, these challenges will typically preclude nonparametric estimation as a feasible strategy

- Thus, we will focus exclusively on parametric approaches to estimation

# Estimation with linear models

- Consider the following set of linear and additive models, where $c^\perp = c - \bar{C}$:

$$E(M_k|c,d) = \beta_{0k} + \beta_{1k}^T c^\perp + \beta_{2k}d \qquad \text{for } k = 1, \ldots, K$$

$$E(Y|c,d,\mathbf{m}) = \gamma_0 + \gamma_1^T c^\perp + \gamma_2 d + \sum_{k=1}^K \gamma_{3k}m_k$$

- Under these models, the natural effects of interest are given by:

$$MNDE(d,d^*) = \gamma_2(d - d^*)$$

$$MNIE(d,d^*) = \left(\sum_{k=1}^K \beta_{2k}\gamma_{3k}\right)(d - d^*)$$

- To compute effect estimates, fit these models by OLS and plug the parameter estimates into the expressions above

# Estimation with linear models

- Now consider the following set of linear models with $D \times M_k$ interactions:

$$E(M_k|c,d) = \beta_{0k} + \beta_{1k}^T c^\perp + \beta_{2k} d \qquad\qquad \text{for } k = 1, \dots, K$$

$$E(Y|c,d,\mathbf{m}) = \gamma_0 + \gamma_1^T c^\perp + \gamma_2 d + \sum_{k=1}^{K} m_k (\gamma_{3k} + \gamma_{4k} d)$$

- Under these models, the natural effects of interest are given by:

$$MNDE(d, d^*) = \left( \gamma_2 + \sum_{k=1}^{K} \gamma_{4k} (\beta_{0k} + \beta_{2k} d^*) \right)(d - d^*)$$

$$MNIE(d, d^*) = \left( \sum_{k=1}^{K} \beta_{2k} (\gamma_{3k} + \gamma_{4k} d) \right)(d - d^*)$$

- To compute effect estimates, fit these models by OLS and plug the parameter estimates into the expressions above

# Estimation with linear models

- Lastly, consider the following set of linear models with covariate interactions:

$$E(M_k|c,d) = \beta_{0k} + \beta_{1k}^T c^\perp + d\left(\beta_{2k} + \beta_{3k}^T c^\perp\right) \qquad \text{for } k = 1,\ldots,K$$

$$E(Y|c,d,\mathbf{m}) = \gamma_0 + \gamma_1^T c^\perp + \gamma_2 d + \sum_{k=1}^K m_k(\gamma_{3k} + \gamma_{4k}d) + c^\perp \sum_{k=1}^K \left(\gamma_{5k}^T d + m_k\left(\gamma_{6k}^T + \gamma_{7k}^T d\right)\right)$$

- Under these models, the natural effects of interest are given by the same expressions as before, provided that the baseline confounders have been mean centered:

$$MNDE(d,d^*) = \left(\gamma_2 + \sum_{k=1}^K \gamma_{4k}(\beta_{0k} + \beta_{2k}d^*)\right)(d - d^*)$$

$$MNIE(d,d^*) = \left(\sum_{k=1}^K \beta_{2k}(\gamma_{3k} + \gamma_{4k}d)\right)(d - d^*)$$

# Summary

- Natural direct and indirect effects through multiple mediators can be estimated using linear models for the mediators and outcome fit to sample data by the method of least squares

- These estimators are consistent provided that the assumptions required for identification are satisfied and provided that all the models used for estimation are correctly specified

- They can easily accommodate exposure-mediator interactions and effect moderation across levels of the baseline confounders

# Limitations

- Models that are linear in the parameters may not perform very well when any of the mediators or the outcome is binary, ordinal, nominal, or a count

  - This approach is best suited for applications in which the mediators and outcome are unbounded and possess equal-interval scaling

  - Nevertheless, there are some situations where a linear model can provide a reasonable approximation for the conditional expected value of a binary, ordinal, or count variable, in which case this approach to estimation remains defensible

- Although this approach easily accommodates exposure-mediator and covariate interactions, it is much more difficult to incorporate interactions among the different mediators, and naïve attempts to do so can lead to uncongenial models

# Estimation via simulation

- Multivariate natural direct and indirect effects can also be estimated using a simulation approach that is implemented with generalized linear models (GLMs)

- The class of GLMs is broad and subsumes normal linear regression as a special case; it also includes a number of nonlinear models, such as logit, probit, and Poisson regression, among others

- This approach to estimation is therefore very general and can be used in a wide variety of different applications (i.e., with continuous, binary, ordinal, nominal, or count variables)

# Estimation via simulation

- The simulation estimator is implemented through a series of steps:

    1. Fit models for each of the mediators

    2. Simulate potential values for each of the mediators

    3. Fit a model for the outcome

    4. Simulate potential outcomes using the simulated values of the mediators

    5. Compute effect estimates using the simulated outcomes

# Estimation via simulation

- Step 1: fit models for each of the mediators

  - Fit a GLM for each mediator, given the baseline confounders, the exposure, and all preceding mediators, denoted by $g_k(M_k|C, D, \mathbf{M}_{k-1})$ where $\mathbf{M}_{k-1} = \{M_1, \dots, M_{k-1}\}$

    - For example:

      $$g_1(M_1|c, d) = Bern\left(p = \text{logit}^{-1}(\beta_{01} + \beta_{11}^T c + \beta_{21}d)\right)$$

      $$g_2(M_2|c, d, m_1) = Normal(\mu = (\beta_{02} + \beta_{12}^T c + \beta_{22}d + \beta_{32}m_1), \sigma^2)$$

  - Let $\hat{g}_k(M_k|C, D, \mathbf{M}_{k-1})$ denote these models with their parameters estimated by maximum likelihood

# Estimation via simulation

- Step 2: simulate potential values for the mediator

  - For every individual in the sample…

    - First, simulate one copy of $M_1(d^*)$ from $\hat{g}_1(M_1|C, d^*)$, and then simulate one copy of $M_1(d)$ from $\hat{g}_1(M_1|C, d)$; let $\widetilde{M}_1(d^*)$ and $\widetilde{M}_1(d)$ denote these simulated values

    - Next, for all $k > 1$ mediators, simulate one copy of $M_k(d^*)$ from $\hat{g}_k\big(M_k|C, d^*, \widetilde{\mathbf{M}}_{k-1}(d^*)\big)$ and one copy of $M_k(d)$ from $\hat{g}_k\big(M_k|C, d, \widetilde{\mathbf{M}}_{k-1}(d)\big)$, where $\widetilde{\mathbf{M}}_{k-1}(d) = \big\{\widetilde{M}_1(d), \dots, \widetilde{M}_{k-1}(d)\big\}$ and $\widetilde{\mathbf{M}}_{k-1}(d^*)$ is defined analogously

    - Repeat these steps $10^3 \leq J \leq 10^4$ times

  - Let $\widetilde{M}_{jk}(d^*)$ and $\widetilde{M}_{jk}(d)$ denote the simulated values for each mediator $k = 1, \dots, K$ and for each simulation $j = 1, 2, \dots, J$, and let $\widetilde{\mathbf{M}}_j(d) = \big\{\widetilde{M}_{j1}(d), \dots, \widetilde{M}_{jK}(d)\big\}$ denote a vector of simulated mediators, with $\widetilde{\mathbf{M}}_j(d^*)$ defined analogously

# Estimation via simulation

- Step 3: fit a model for the outcome

  - Fit a GLM for the outcome given the baseline confounders, the exposure, and the vector of mediators, denoted by $h(Y|C, D, \mathbf{M})$, where $\mathbf{M} = \{M_1, \ldots, M_K\}$

    - For example:

$$h(Y|c, d, \mathbf{m}) = Pois\left(\lambda = \exp\left(\gamma_0 + \gamma_1^T c^\perp + \gamma_2 d + \sum_{k=1}^{K} m_k (\gamma_{3k} + \gamma_{4k} d)\right)\right)$$

  - Let $\hat{h}(Y|C, D, \mathbf{M})$ denote this model with its parameters estimated by maximum likelihood

# Estimation via simulation

- Step 4: simulate potential outcomes

  - For every sample member and each simulated vector of mediators…

    - simulate one copy of $Y(d, \mathbf{M}(d))$ from $\hat{h}(Y|C, d, \widetilde{\mathbf{M}}_j(d))$ and then…

    - simulate one copy of $Y(d^*, \mathbf{M}(d^*))$ from $\hat{h}(Y|C, d^*, \widetilde{\mathbf{M}}_j(d^*))$ and then…

    - simulate one copy of $Y(d, \mathbf{M}(d^*))$ from $\hat{h}(Y|C, d, \widetilde{\mathbf{M}}_j(d^*))$

  - Let $\tilde{Y}_j(d, \mathbf{M}(d))$, $\tilde{Y}_j(d^*, \mathbf{M}(d^*))$, and $\tilde{Y}_j(d, \mathbf{M}(d^*))$ denote the simulated values of the outcome for each simulation $j = 1, 2, \ldots, J$

# Estimation via simulation

- Step 4: compute effect estimates

  - Average the differences between simulated outcomes over simulations and over sample members as follows…

  $$\widehat{MNDE}(d, d^*) = \frac{1}{nJ}\sum\sum_j\left[\tilde{Y}_j\big(d, \mathbf{M}(d^*)\big) - \tilde{Y}_j\big(d^*, \mathbf{M}(d^*)\big)\right]$$

  $$\widehat{MNIE}(d, d^*) = \frac{1}{nJ}\sum\sum_j\left[\tilde{Y}_j\big(d, \mathbf{M}(d)\big) - \tilde{Y}_j\big(d, \mathbf{M}(d^*)\big)\right]$$

  $$\widehat{ATE}(d, d^*) = \frac{1}{nJ}\sum\sum_j\left[\tilde{Y}_j\big(d, \mathbf{M}(d)\big) - \tilde{Y}_j\big(d^*, \mathbf{M}(d^*)\big)\right]$$

# Model specification

- This approach can easily accommodate exposure-mediator interactions, mediator-mediator interactions, covariate interactions, and nonlinear terms, as well as many different link functions and distribution models

- The steps outlined previously proceed exactly the same, regardless of the particular form of the GLMs used for the mediators and outcome

# Summary

- Multivariate natural direct and indirect effects can be estimated via simulation with a broad class of GLMs fit to sample data by the method of maximum likelihood

- These estimators are consistent provided that the assumptions required for identification are satisfied and provided that all the models used for estimation are correctly specified

- Limitations

  - The method requires correctly specified models for all the mediators and the outcome, which may be difficult to achieve in practice, especially in applications with many mediators

# Estimation via weighting

- In contrast to linear models and the simulation approach, which both require models for all the mediators and the outcome, weighting estimators are implemented only with models for the exposure

- These models are used to construct a set of weights that transform the empirical distribution of the sample data in ways that emulate different hypothetical experiments

- The effects of interest are estimated by comparing the mean of the outcome across differently weighted samples

# Estimation via weighting

- The weighting estimator is implemented through a series of steps:

  1. Fit two different models for the exposure

  2. Compute predicted probabilities of exposure from each model

  3. Use the exposure probabilities to construct a set of weights

  4. Compute effect estimates by comparing weighted means of the outcome

# Estimation via weighting

- Step 1: fit models for the exposure

  - Fit a GLM for the exposure given the baseline confounders, denoted by $f(D|C)$

  - Next, fit another GLM for the exposure given the baseline confounders and the vector of mediators, denoted by $s(D|C, \mathbf{M})$, where $\mathbf{M} = \{M_1, \dots, M_K\}$

    - If, for example, the exposure is binary, then $f(D|C)$ and $s(D|C, \mathbf{M})$ might be logit or probit models

  - Let $\hat{f}(D|C)$ and $\hat{s}(D|C, \mathbf{M})$ denote these models with their parameters estimated by maximum likelihood

# Estimation via weighting

- Step 2: compute predicted probabilities of exposure

  - For each sample member, use $\hat{f}(D|C)$ to predict…

    - the probability of exposure to $d$ given their baseline confounders, denoted by $\hat{P}(d|C)$

    - the probability of exposure to $d^*$ given their baseline confounders, denoted by $\hat{P}(d^*|C)$

  - Next, for each sample member, use $\hat{s}(D|C, \mathbf{M})$ to predict…

    - the probability of exposure to $d$ given their baseline confounders and values on all the mediators, denoted by $\hat{P}(d|C, \mathbf{M})$

    - the probability of exposure to $d^*$ given their baseline confounders and values on all the mediators, denoted by $\hat{P}(d^*|C, \mathbf{M})$

# Estimation via weighting

- Step 3: construct IPWs

  - Among sample members with $D = d^*$, compute...

    - $\widehat{wm}_1 = \dfrac{1}{\hat{P}(d^*|C)}$

  - Among sample members with $D = d$, compute...

    - $\widehat{wm}_2 = \dfrac{1}{\hat{P}(d|C)}$

    - $\widehat{wm}_3 = \dfrac{\hat{P}(d^*|C, \mathbf{M})}{\hat{P}(d|C, \mathbf{M})\hat{P}(d^*|C)}$

# Estimation via weighting

- Step 4: compute effect estimates

    - Compute differences between weighted means of the observed outcome as follows…

$$\widehat{MNDE}(d, d^*) = \frac{\sum I(D=d)\widehat{wm}_3 Y}{\sum I(D=d)\widehat{wm}_3} - \frac{\sum I(D=d^*)\widehat{wm}_1 Y}{\sum I(D=d^*)\widehat{wm}_1}$$

$$\widehat{MNIE}(d, d^*) = \frac{\sum I(D=d)\widehat{wm}_2 Y}{\sum I(D=d)\widehat{wm}_2} - \frac{\sum I(D=d)\widehat{wm}_3 Y}{\sum I(D=d)\widehat{wm}_3}$$

$$\widehat{ATE}(d, d^*) = \frac{\sum I(D=d)\widehat{wm}_2 Y}{\sum I(D=d)\widehat{wm}_2} - \frac{\sum I(D=d^*)\widehat{wm}_1 Y}{\sum I(D=d^*)\widehat{wm}_1}$$

# Stabilized and censored weights

- Stabilized versions of the inverse probability weights can be expressed as follows:

  - $\widehat{swm}_1 = \dfrac{\hat{P}(d^*)}{\hat{P}(d^*|C)}$

  - $\widehat{swm}_2 = \dfrac{\hat{P}(d)}{\hat{P}(d|C)}$

  - $\widehat{swm}_3 = \dfrac{\hat{P}(d^*|C, \mathbf{M})\hat{P}(d^*)}{\hat{P}(d|C, \mathbf{M})\hat{P}(d^*|C)}$

- The performance of these weights can usually be improved even further by censoring their extreme values—for example, at the 1st and 99th percentiles

# Summary

- Multivariate natural direct and indirect effects can be estimated via weighting with two different GLMs for the probability of exposure

- These estimators are consistent provided that the assumptions required for identification are satisfied and provided that the models used for the exposure are correctly specified

- Limitations

  - Difficult to use and often unstable with continuous or many valued exposures

  - Highly sensitive to model misspecification

# Regression imputation

- Multivariate natural direct and indirect effects can also be estimated using a regression imputation approach

- Unlike the other approaches we've considered, regression imputation does not require models for the mediators or for the exposure; rather, it only requires a series of models for the outcome

# Regression imputation

- Regression imputation is implemented through a series of steps:

    1. Fit a model for the outcome given the exposure and baseline confounders
        - Impute outcomes from this model under $D = d$ and $D = d^*$

    2. Fit another model for the outcome given the exposure, confounders, and mediators
        - Impute outcomes from this model under $D = d$

    3. Fit a third model for the imputed outcomes from the prior step
        - Impute outcomes from this model under $D = d^*$

    4. Compute effect estimates using the different imputed outcomes

# Regression imputation

- Step 1: Fit a model for the outcome and construct imputations

    - Fit a model for the outcome given the baseline confounders and the exposure, denoted by $q(Y|C, D)$

        - Let $\hat{q}(Y|C, D)$ denote this model with its parameters estimated by least squares or maximum likelihood

    - Impute potential outcomes under $d^*$ by setting $D = d^*$ for all sample members and computing predicted values, given by $\hat{Y}(d^*) = \hat{q}(C, d^*)$

    - Impute potential outcomes under $d$ by setting $D = d$ for all sample members and computing predicted values, given by $\hat{Y}(d) = \hat{q}(C, d)$

# Regression imputation

- Step 2: Fit another model for the outcome and construct imputations

    - Fit a model for the outcome given the baseline confounders, the exposure, and the vector of mediators, denoted by $h(Y|C, D, \mathbf{M})$

        - Let $\hat{h}(Y|C, D, \mathbf{M})$ denote this model with its parameters estimated by least squares or maximum likelihood

    - Impute potential outcomes under $d$ and $\mathbf{M}(D)$ by setting $D = d$ for all sample members and computing predicted values, given by $\hat{Y}\big(d, \mathbf{M}(D)\big) = \hat{h}(C, d, \mathbf{M})$

# Regression imputation

- Step 3: fit a model for the imputed outcomes and construct imputations

    - Fit a model for the imputed outcomes constructed in the previous step given the baseline confounders and exposure, denoted by $\tau\big(\hat{Y}(d, \mathbf{M}(D))\big|C, D\big)$

        - Let $\hat{\tau}\big(\hat{Y}(d, \mathbf{M}(D))\big|C, D\big)$ denote this model with its parameters estimated by least squares or maximum likelihood

    - Impute potential outcomes under $d$ and $\mathbf{M}(d^*)$ by setting $D = d^*$ for all sample members and computing predicted values, given by $\hat{Y}\big(d, \mathbf{M}(d^*)\big) = \hat{\tau}\big(\hat{Y}(d, \mathbf{M}(D))\big|C, d^*\big)$

# Regression imputation

- Step 4: compute effect estimates

  - Compute differences between means of the imputed outcomes as follows…

$$\widehat{MNDE}(d, d^*) = \frac{1}{n}\sum\left[\hat{Y}(d, \mathbf{M}(d^*)) - \hat{Y}(d^*)\right]$$

$$\widehat{MNIE}(d, d^*) = \frac{1}{n}\sum\left[\hat{Y}(d) - \hat{Y}(d, \mathbf{M}(d^*))\right]$$

$$\widehat{ATE}(d, d^*) = \frac{1}{n}\sum\left[\hat{Y}(d) - \hat{Y}(d^*)\right]$$

# Summary

- Multivariate natural direct and indirect effects can be estimated via regression imputation with a series of linear models for the outcome

- These estimators are consistent provided that the assumptions required for identification are satisfied and provided that the models used for estimation are correctly specified

- Like the weighting approach, regression imputation is especially useful in analyses of multiple mediators because it obviates the need to correctly specify and fit a model for each mediator

# Example: NLSY79

- 1979 National Longitudinal Study of Youth

  - Exposure ($D$)
    - sample member attended college before age 22

  - Outcome ($Y$):
    - standardized scores on the CES-D at age 40

  - Covariates ($C$):
    - race, gender, parental education, occupation, and income, household size, AFQT scores

  - Potential mediators (**M**)
    - unemployment between age 35-40 ($M_1$)
    - household income between age 35-40 ($M_2$)

# Example: NLSY79

- Many studies have documented that going to college seems to reduce the likelihood of becoming depressed later in life—but how does this effect come about?

- One possibility is that a more advanced education reduces depression by increasing the labor market prospects of adults, boosting both their employment and wages

  - Do unemployment and income jointly mediate the effect of college attendance on depression?

# Example: NLSY79

- Using linear models, compute estimates for the *MNDE* and *MNIE* of education on depression operating through income and unemployment

```
1   ### wk 7 nlsy tutorial ###
2   rm(list=ls())
3
4   ## load/install libraries ##
5   packages<-c("dplyr", "tidyr", "foreign", "foreach", "doParallel", "doRNG", "devtools")
6   install.packages(packages)
7
8 ▾ for (package.i in packages) {
9     suppressPackageStartupMessages(library(package.i, character.only=TRUE))
10 ▴   }
11
12  ## load data ##
13  datadir <- "C:/Users/Geoffrey Wodtke/Dropbox/D/courses/2024-25_UOFCHICAGO/SOCI_40258_CAUSAL_MEDIATION/data/"
14  nlsy <- read.dta(paste(datadir, "nlsy79.dta", sep=""))
15
16  Y <- "std_cesd_age40"
17  D <- "att22"
18  M1 <- "ever_unemp_age3539"
19  M2 <- "log_faminc_adj_age3539"
20  C <- c("female", "black", "hispan", "paredu", "parprof", "parinc_prank", "famsize", "afqt3")
21
22  nlsy <- nlsy[complete.cases(nlsy[,c(C,D,M1,M2,"cesd_age40")]),] |>
23    mutate(std_cesd_age40 = (cesd_age40 - mean(cesd_age40)) / sd(cesd_age40))|
```

# Example: NLSY79

- Using linear models, compute estimates for the *MNDE* and *MNIE* of education on depression operating through income and unemployment

```
25  ## compute estimates w/ linear models ##
26
27  #load R functions
28  source("https://raw.githubusercontent.com/causalMedAnalysis/causalMedR/refs/heads/main/utils.R")
29  source("https://raw.githubusercontent.com/causalMedAnalysis/causalMedR/refs/heads/main/linmed.R")
30
31  #compute estimates
32  lin_est <- linmed(data = nlsy, D = D, M = c(M1, M2), Y = Y, C = C,
33    interaction_DM = TRUE, interaction_DC = TRUE, interaction_MC = TRUE,
34    boot = TRUE, boot_reps = 2000, boot_seed = 60637, boot_parallel = TRUE)
```

# Example: NLSY79

- Using linear models, compute estimates for the *MNDE* and *MNIE* of education on depression operating through income and unemployment

```
36  lin_output <- data.frame(
37    param = c("ATE(1,0)", "MNDE(1,0)", "MNIE(1,0)"),
38    est = c(lin_est$ATE, lin_est$NDE, lin_est$NIE),
39    ci_lo = c(lin_est$ci_ATE[1], lin_est$ci_NDE[1], lin_est$ci_NIE[1]),
40    ci_hi = c(lin_est$ci_ATE[2], lin_est$ci_NDE[2], lin_est$ci_NIE[2]),
41    pval = c(lin_est$pvalue_ATE, lin_est$pvalue_NDE, lin_est$pvalue_NIE)) |>
42    mutate(across(.cols = !param, .fns = \(x) round(x, 3)))
43
44  print(lin_output)
```

```
> print(lin_output)
        param    est   ci_lo  ci_hi  pval
1   ATE(1,0) -0.113 -0.206 -0.014 0.020
2 MNDE(1,0) -0.035 -0.134  0.073 0.519
3 MNIE(1,0) -0.078 -0.121 -0.042 0.000
```

# Example: NLSY79

- Using the simulation approach, compute estimates for the *MNDE* and *MNIE* of education on depression through income and unemployment

```
48  #load R functions
49  source("https://raw.githubusercontent.com/causalMedAnalysis/causalMedR/refs/heads/main/medsim.R")
50
51  #specify models for M1 (logit), M2 (normal linear), and Y (normal linear)
52  formula_M1 <- paste(M1, "~",
53    paste(paste(c(D, C), collapse = " + "), "+",
54    paste(D, C, sep = ":", collapse = " + ")))
55
56  formula_M2 <- paste(M2, "~",
57    paste(paste(paste(paste(c(D, M1, C), collapse = " + "), "+",
58    paste(D, M1, sep = ":", collapse = " + ")), "+",
59    paste(D, C, sep = ":", collapse = " + ")), "+",
60    paste(M1, C, sep = ":", collapse = " + ")))
61
62  formula_Y <- paste(Y, "~",
63    paste(paste(paste(paste(paste(paste(c(D, M1, M2, C), collapse = " + "), "+",
64    paste(D, M1, sep = ":", collapse = " + ")), "+",
65    paste(D, M2, sep = ":", collapse = " + ")), "+",
66    paste(D, C, sep = ":", collapse = " + ")), "+",
67    paste(M1, C, sep = ":", collapse = " + ")), "+",
68    paste(M2, C, sep = ":", collapse = " + ")))
69
```

# Example: NLSY79

- Using the simulation approach, compute estimates for the *MNDE* and *MNIE* of education on depression through income and unemployment

```
70  specs <- list(
71     list(func = "glm", formula = as.formula(formula_M1), args = list(family = "binomial")),
72     list(func = "lm", formula = as.formula(formula_M2)),
73     list(func = "lm", formula = as.formula(formula_Y)))
74
75  #compute estimates
76  sim_est <- medsim(data = nlsy, num_sim = 1000, treatment = D, intv_med = NULL,
77     model_spec = specs, seed = 60637, boot = TRUE, reps = 2000)
78
```

# Example: NLSY79

- Using the simulation approach, compute estimates for the *MNDE* and *MNIE* of education on depression through income and unemployment

```
79  sim_output <- data.frame(
80    param = c("ATE(1,0)", "MNDE(1,0)", "MNIE(1,0)"),
81    est = c(sim_est$point.est[1], sim_est$point.est[2], sim_est$point.est[3]),
82    ci_lo = c(sim_est$ll.95ci[1], sim_est$ll.95ci[2], sim_est$ll.95ci[3]),
83    ci_hi = c(sim_est$ul.95ci[1], sim_est$ul.95ci[2], sim_est$ul.95ci[3]),
84    pval = c(sim_est$pval[1], sim_est$pval[2], sim_est$pval[3])) |>
85    mutate(across(.cols = !param, .fns = \(x) round(x, 3)))
86
87  print(sim_output)
```

```
> print(sim_output)
       param    est  ci_lo  ci_hi  pval
1   ATE(1,0) -0.119 -0.215 -0.020 0.012
2 MNDE(1,0) -0.036 -0.135  0.073 0.523
3 MNIE(1,0) -0.084 -0.133 -0.044 0.000
```

# Example: NLSY79

- Using inverse probability weights, compute estimates for the *MNDE* and *MNIE* of education on depression through income and unemployment

```
89   ## compute estimates w/ inverse probability weighting ##
90
91   #load R functions
92   source("https://raw.githubusercontent.com/causalMedAnalysis/causalMedR/refs/heads/main/ipwmed.R")
93
94   #specify models for D
95   f_of_D_giv_C <- paste(D, "~", paste(C, collapse = " + "))
96   s_of_D_giv_CM1M2 <- paste(D, "~", paste(c(M1, M2, C), collapse = " + "))
97
98   #compute estimates
99   ipw_est <- ipwmed(data = nlsy, D = D, M = c(M1, M2), Y = Y,
100     formula1_string = f_of_D_giv_C, formula2_string = s_of_D_giv_CM1M2,
101     stabilize = TRUE, censor = TRUE,
102     boot = TRUE, boot_reps = 2000, boot_seed = 60637, boot_parallel = TRUE)
```

# Example: NLSY79

- Using inverse probability weights, compute estimates for the *MNDE* and *MNIE* of education on depression through income and unemployment

```
104  ipw_output <- data.frame(
105    param = c("ATE(1,0)", "MNDE(1,0)", "MNIE(1,0)"),
106    est = c(ipw_est$ATE, ipw_est$NDE, ipw_est$NIE),
107    ci_lo = c(ipw_est$ci_ATE[1], ipw_est$ci_NDE[1], ipw_est$ci_NIE[1]),
108    ci_hi = c(ipw_est$ci_ATE[2], ipw_est$ci_NDE[2], ipw_est$ci_NIE[2]),
109    pval = c(ipw_est$pvalue_ATE, ipw_est$pvalue_NDE, ipw_est$pvalue_NIE)) |>
110    mutate(across(.cols = !param, .fns = \(x) round(x, 3)))
111
112  print(ipw_output)
```

```
> print(ipw_output)
        param     est  ci_lo  ci_hi  pval
1   ATE(1,0) -0.167 -0.264 -0.053 0.003
2  MNDE(1,0) -0.100 -0.227  0.055 0.188
3  MNIE(1,0) -0.068 -0.141 -0.009 0.025
```

# Example: NLSY79

- Using regression imputation, compute estimates for the *MNDE* and *MNIE* of education on depression operating through income and unemployment

```
114  ## compute estimates w/ regression imputation ##
115
116  #define regression imputation function
117 ▾ impmed <- function(data) {
118
119    df <- data
120
121    Ymodel_CD <- lm(std_cesd_age40 ~ att22 * (female + black + hispan +
122      paredu + parprof + parinc_prank + famsize + afqt3), data=df)
123
124    idata <- df
125
126    idata$att22 <- 0
127
128    Y0hat <- predict(Ymodel_CD, newdata=idata, type="response")
129
130    idata$att22 <- 1
131
132    Y1hat <- predict(Ymodel_CD, newdata=idata, type="response")
```

# Example: NLSY79

- Using regression imputation, compute estimates for the *MNDE* and *MNIE* of education on depression operating through income and unemployment

```
134   Ymodel_CDM <- lm(std_cesd_age40 ~
135     att22 * (female + black + hispan + paredu + parprof + parinc_prank +
136     famsize + afqt3 + log_faminc_adj_age3539 + ever_unemp_age3539), data=df)
137
138   idata <- df
139
140   idata$att22 <- 1
141
142   df$Y1MDhat <- predict(Ymodel_CDM, newdata=idata, type="response")
143
144   YhatModel_CD <- lm(Y1MDhat ~ att22 * (female + black + hispan +
145     paredu + parprof + parinc_prank + famsize + afqt3), data=df)
146
147   idata <- df
148
149   idata$att22 <- 0
150
151   Y1M0hat <- predict(YhatModel_CD, newdata=idata, type="response")
```

```
154   MNDE <- mean(Y1M0hat) - mean(Y0hat)
155   MNIE <- mean(Y1hat) - mean(Y1M0hat)
156   ATE <- MNDE + MNIE
157
158   point.est <- list(ATE, MNDE, MNIE)
159
160   return(point.est)
161 ▲ }
```

# Example: NLSY79

- Using regression imputation, compute estimates for the *MNDE* and *MNIE* of education on depression operating through income and unemployment

```
163  #compute point estimates
164  impmed.est <- impmed(nlsy)
165  impmed.est <- matrix(unlist(impmed.est), ncol=3, byrow=TRUE)
166
167  #compute bootstrap estimates
168  ncores <- detectCores()-2
169  my.cluster <- parallel::makeCluster(ncores, type="PSOCK")
170  doParallel::registerDoParallel(cl=my.cluster)
171  clusterExport(cl=my.cluster, list("impmed"), envir=environment())
172  registerDoRNG(60637)
173
174 ▾ impmed.boot <- foreach(i=1:2000, .combine=cbind) %dopar% {
175
176      boot.data <- nlsy[sample(nrow(nlsy), nrow(nlsy), replace=TRUE),]
177
178      boot.est <- impmed(data=boot.data)
179                                                          doParallel::regist
180      return(boot.est)
181 ▴ }
182
183  stopCluster(my.cluster)
184  rm(my.cluster)
185
186  impmed.boot <- matrix(unlist(impmed.boot), ncol=3, byrow=TRUE)
```

# Example: NLSY79

- Using regression imputation, compute estimates for the *MNDE* and *MNIE* of education on depression operating through income and unemployment

```
188  #collate estimates
189  impmed.output <- data.frame(
190    param = c("ATE(1,0)", "MNDE(1,0)", "MNIE(1,0)"),
191    est = impmed.est[1:3],
192    ci_lo = apply(impmed.boot, 2, function(x) quantile(x, prob=0.025)),
193    ci_hi = apply(impmed.boot, 2, function(x) quantile(x, prob=0.975))) %>%
194      mutate(across(c(est, ci_lo, ci_hi), ~round(.x, digits = 3)))
195
196  print(impmed.output)
```

```
> print(impmed.output)
         param     est  ci_lo  ci_hi
1    ATE(1,0) -0.119 -0.212 -0.023
2  MNDE(1,0) -0.062 -0.160  0.040
3  MNIE(1,0) -0.056 -0.093 -0.027
```