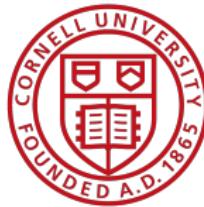


Causal Mediation in Natural Experiments

Senan Hogan-Hennessy
Economics Department, Cornell University
seh325@cornell.edu

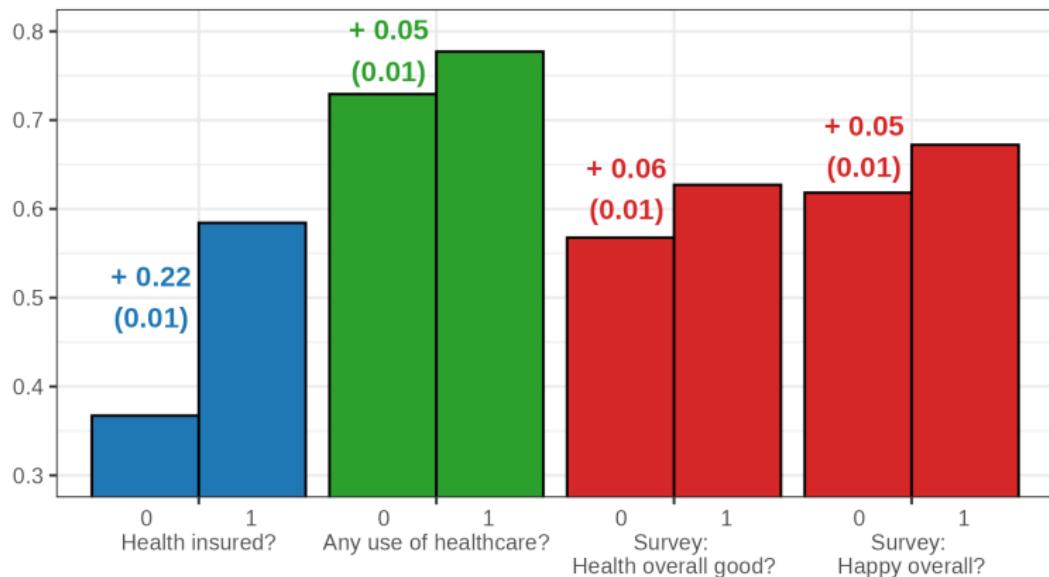


Cornell Placement Week
29 September 2025

Intro: Oregon Health Insurance Experiment

In 2008, Oregon gave access to socialised health insurance by wait-list lottery (Finkelstein et al, 2012).

Mean Outcome, for each $z' = 0, 1$.



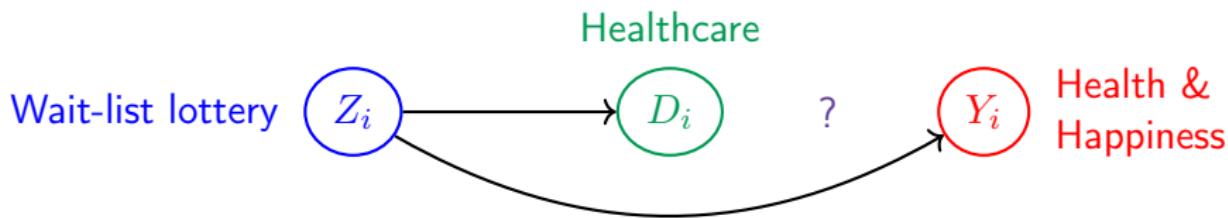
Applied practice:

⇒ Suggestive evidence for healthcare as mechanism for wait-list lottery. . . .

Intro: Oregon Health Insurance Experiment

In 2008, Oregon gave access to socialised health insurance by wait-list lottery (Finkelstein et al, 2012).

Figure: Model for Suggestive Evidence of a Mechanism.



Inconsistencies in suggestive evidence of mechanisms:

- Is $D_i \rightarrow Y_i$ small, large, or even nonexistent?
 - Where else do we accept assumed causal effects without evidence?

Introduction

Causal Mediation (CM) is an alternative framework to studying mechanisms, with clear identification and assumptions required.

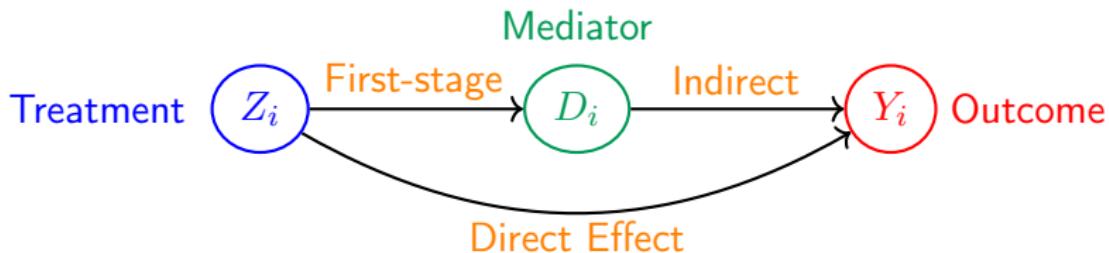
- ① Problems with conventional approach to CM in observational settings.
[Negative result]
 - ② Recovering valid CM effects, via MTE + control function modelling.
[Positive result]
-

New insights from intersection of two fields:

- **CM.**
Imai Keele Yamamoto (2010), Frölich Huber (2017), Deuchert Huber Schelker (2019), Huber (2020), Kwon Roth (2024).
- **Labour theory, Selection-into-treatment, MTEs.**
Roy (1951), Heckman (1979), Heckman Honoré (1990), Vycatil (2002), Heckman Vycatil (2005), Brinch Mogstad Wiswall (2017), Kline Walters (2019).

Introduction – CM

Consider ignorable treatment $Z_i = 0, 1$, binary mediator $D_i = 0, 1$, and continuous outcome Y_i .



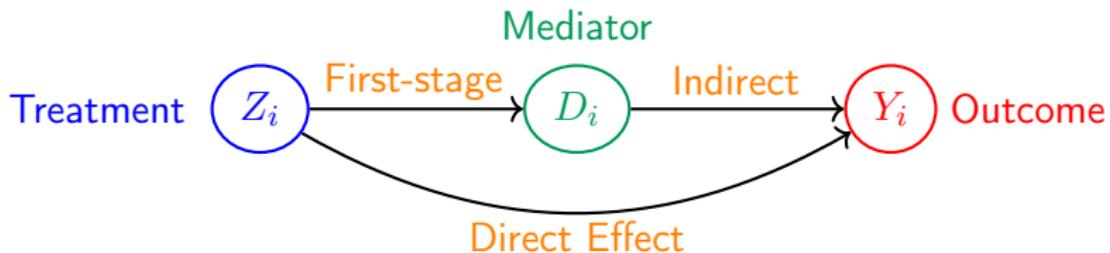
Assumption: Mediator Ignorability (MI, Imai Keele Yamamoto 2010)
mediator D_i is also ignorable, conditional on X_i and Z_i realisation.

Average Direct Effect (ADE) and Average Indirect Effect (AIE) are identified by two-stage regression

- ADE is causal effect $Z_i \rightarrow Y_i$, blocking the indirect D_i path
- AIE is causal effect of $D_i(Z_i) \rightarrow Y_i$, blocking the direct Z_i path.

Introduction – CM

Consider ignorable treatment $Z_i = 0, 1$, binary mediator $D_i = 0, 1$, and continuous outcome Y_i .



Assumption: **Mediator Ignorability** (MI, Imai Keele Yamamoto 2010)
mediator D_i is also ignorable, conditional on X_i and Z_i realisation.

Average Direct Effect (ADE) and Average Indirect Effect (AIE) are identified by two-stage regression

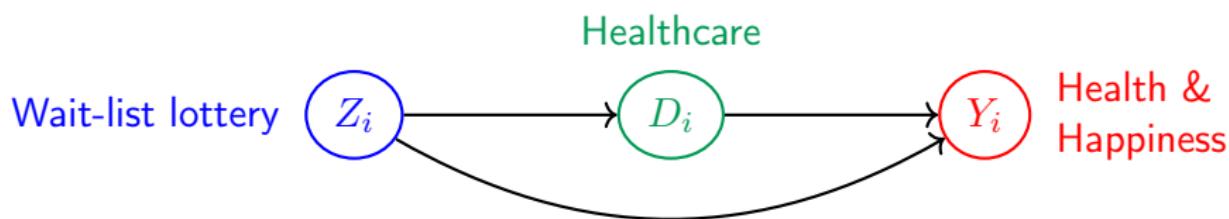
- ADE is causal effect $Z_i \rightarrow Y_i$, blocking the indirect D_i path
- AIE is causal effect of $D_i(Z_i) \rightarrow Y_i$, blocking the direct Z_i path.

1. Selection Bias

Assumption: Mediator ignorability (MI, Imai Keele Yamamoto 2010)
 mediator D_i is *also* ignorable, conditional on X_i, Z_i realisation

Would this assumption hold true in settings economists study?

E.g., Oregon Health Insurance Experiment.



- ① Treatment is as-good-as random (2008 Oregon wait-list lottery).
 - ② Healthcare is quasi-random, conditional on lottery realisation Z_i and demographic controls X_i .

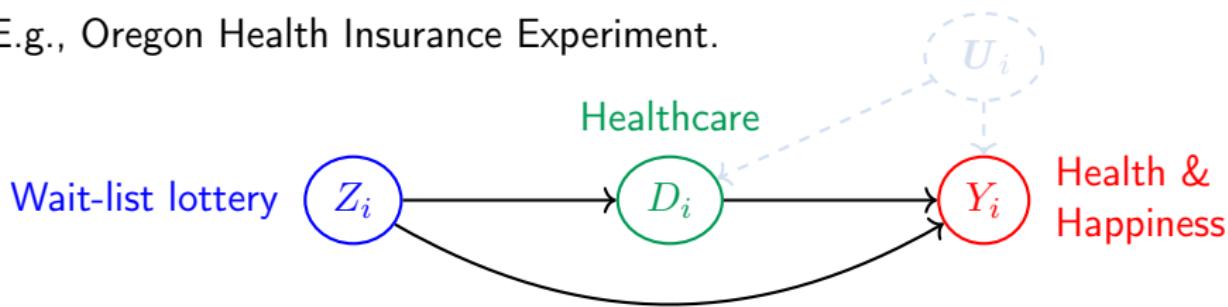
1. Selection Bias

Assumption: Mediator ignorability (MI, Imai Keele Yamamoto 2010)

mediator D_i is **also ignorable**, conditional on X_i , Z_i realisation

Would this assumption hold true in settings economists study?

E.g., Oregon Health Insurance Experiment.



Theorem: If choice to attend healthcare is unconstrained, based on costs and benefits (Roy model) and demographics do not explain all benefits \Rightarrow MI does not hold, there is unobserved confounding.

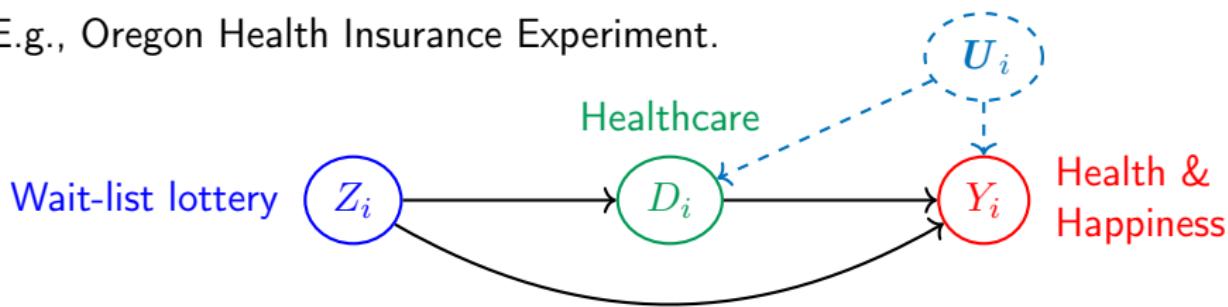
1. Selection Bias

Assumption: Mediator ignorability (MI, Imai Keele Yamamoto 2010)

mediator D_i is *also* ignorable, conditional on X_i , Z_i realisation

Would this assumption hold true in settings economists study?

E.g., Oregon Health Insurance Experiment.



Theorem: If choice to attend healthcare is unconstrained, based on costs and benefits (Roy model) and demographics do not explain all benefits \Rightarrow MI does not hold, there is unobserved confounding.

1. Selection Bias

In an observational setting, need an additional credible research design for **Mediator Ignorability (MI)** to be credible.

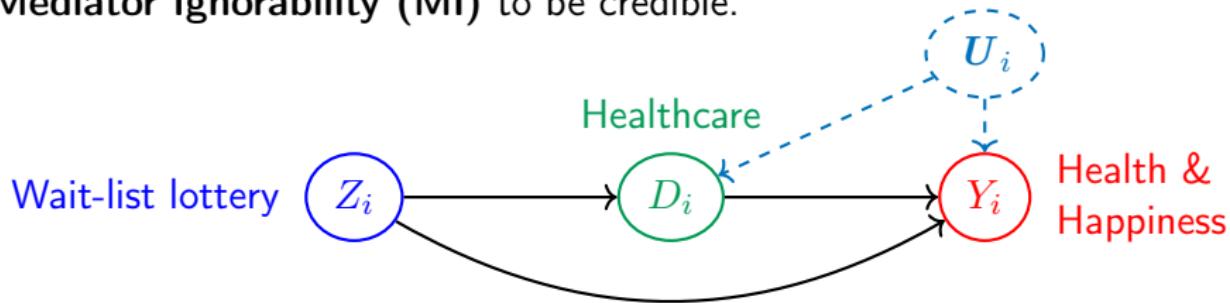
(a) Cells in a lab → MI believable.

(b) People choosing healthcare → MI not.



1. Selection Bias

In an observational setting, need an additional credible research design for **Mediator Ignorability (MI)** to be credible.



If not, then CM effects are contaminated by bias terms, similar to classical selection bias (e.g., Heckman Ichimura Smith Todd 1998).

- ADE: $CM \text{ Estimand} = ADE + (Selection \text{ Bias} + Group \text{ difference bias})$
- AIE: $CM \text{ Estimand} = AIE + (Selection \text{ Bias} + Group \text{ difference bias})$

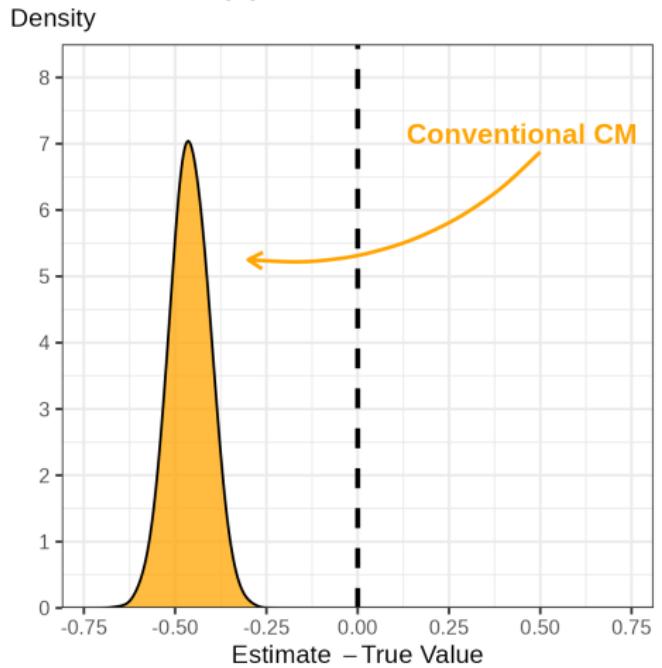
► ADE biases

► AIE biases

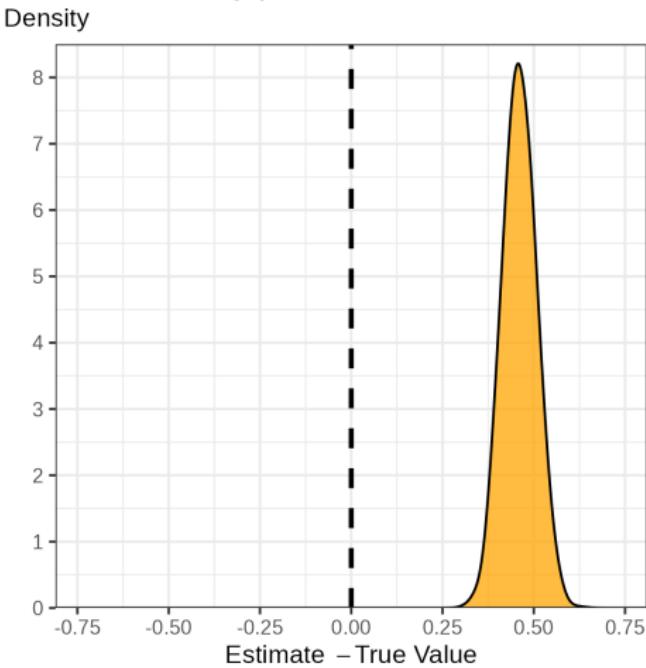
1. Selection Bias

In a simulation with Roy selection-into- D_i , CM estimates are biased.

(a) $\widehat{ADE} - ADE$.



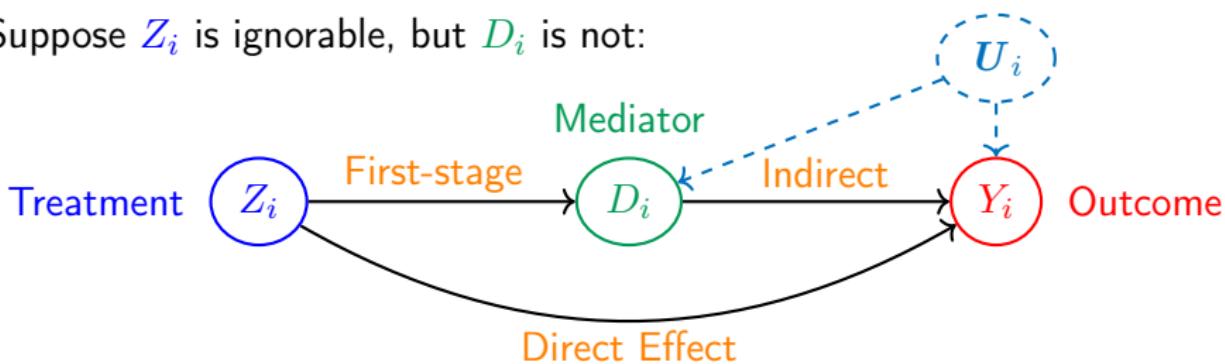
(b) $\widehat{AIE} - AIE$.



2. CM with Selection

Conventional CM methods do not identify ADE + AIE in a natural experiment setting, but can we build a credible structural model?

Suppose Z_i is ignorable, but D_i is not:



- ① Average first-stage, $Z_i \rightarrow D_i$, is identified
 - ② Average second-stage, $Z_i, D_i \rightarrow Y_i$, is not — represented by U_i .

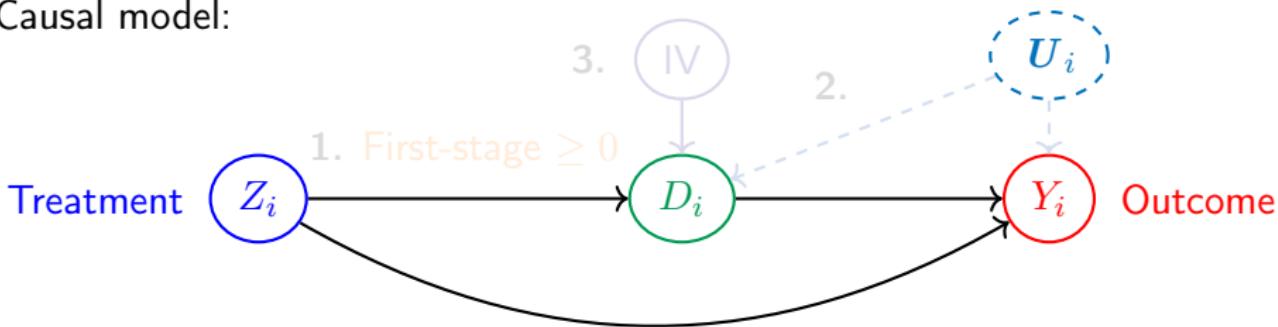
Intuition: model U_i via mediator MTE to identify ADE + AIE.

2. CM with Selection — Identification

MTE assumptions:

- ① Mediator monotonicity
- ② Selection on mediator benefits
- ③ IV for mediator take-up cost.

Causal model:



Proposition: Under MTE assumptions, the mediator MTE is identified.

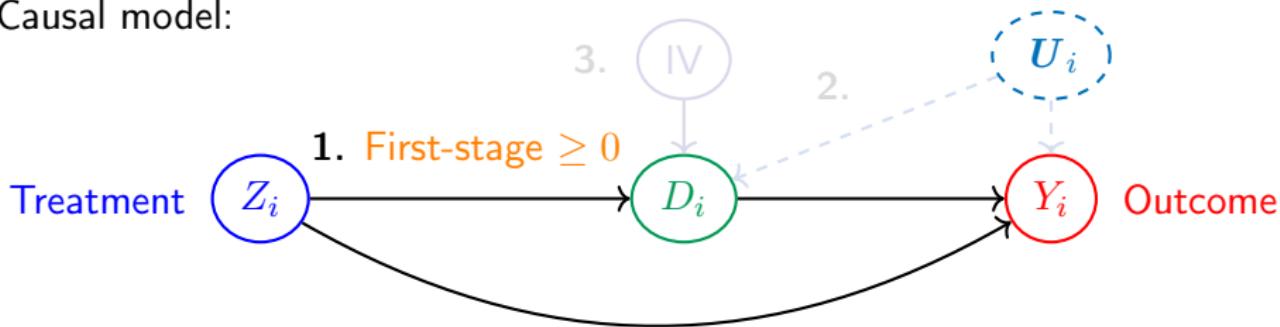
Theorem: Mediation second-stage effects, $Z_i, D_i \rightarrow Y_i$, are identified by the MTE associated Control Functions (CFs).

2. CM with Selection — Identification

MTE assumptions:

- ① Mediator monotonicity
- ② Selection on mediator benefits
- ③ IV for mediator take-up cost.

Causal model:



Proposition: Under MTE assumptions, the mediator MTE is identified.

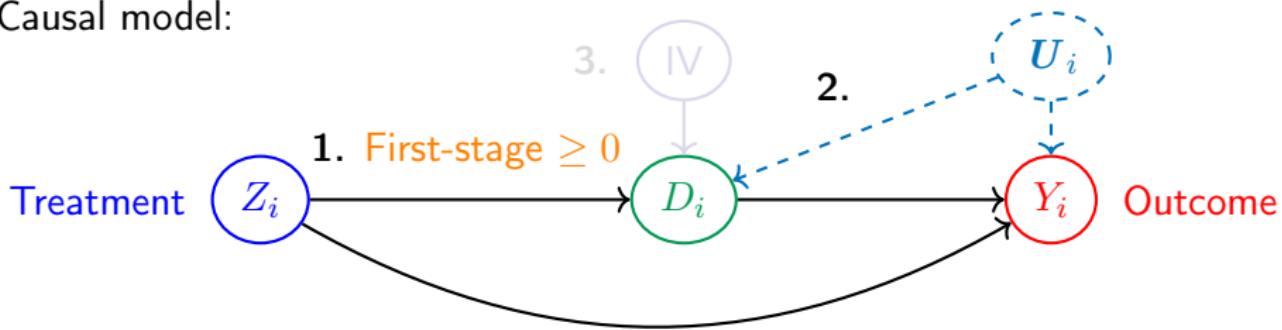
Theorem: Mediation second-stage effects, $Z_i, D_i \rightarrow Y_i$, are identified by the MTE associated Control Functions (CFs).

2. CM with Selection — Identification

MTE assumptions:

- ① Mediator monotonicity
- ② Selection on mediator benefits
- ③ IV for mediator take-up cost.

Causal model:



Proposition: Under MTE assumptions, the mediator MTE is identified.

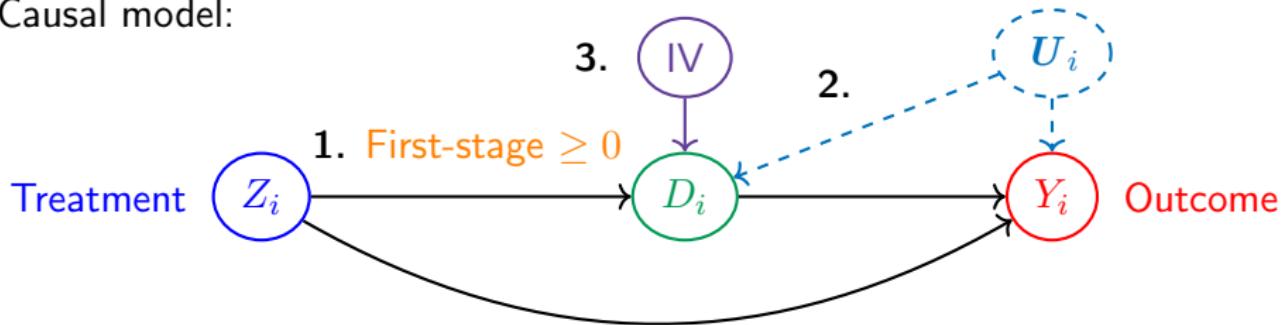
Theorem: Mediation second-stage effects, $Z_i, D_i \rightarrow Y_i$, are identified by the MTE associated Control Functions (CFs).

2. CM with Selection — Identification

MTE assumptions:

- ① Mediator monotonicity
- ② Selection on mediator benefits
- ③ IV for mediator take-up cost.

Causal model:



Proposition: Under MTE assumptions, the mediator MTE is identified.

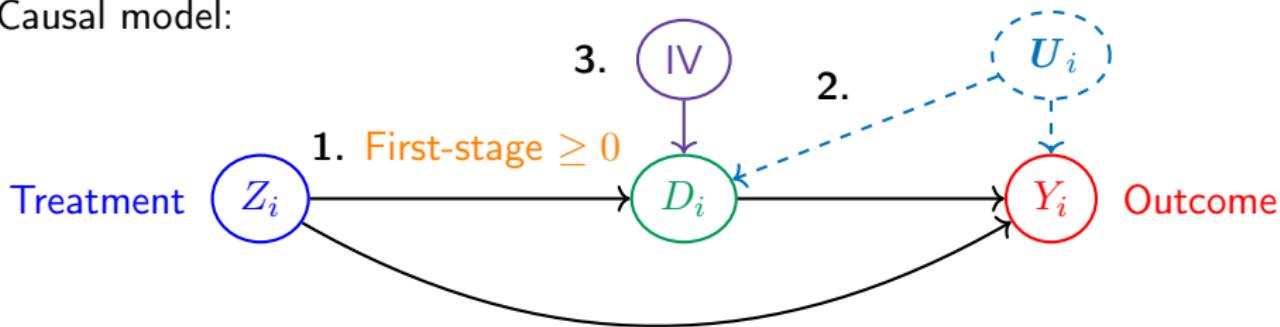
Theorem: Mediation second-stage effects, $Z_i, D_i \rightarrow Y_i$, are identified by the MTE associated Control Functions (CFs).

2. CM with Selection — Identification

MTE assumptions:

- ① Mediator monotonicity
 - ② Selection on mediator benefits
 - ③ IV for mediator take-up cost.

Causal model:



Proposition: Under MTE assumptions, the mediator MTE is identified.

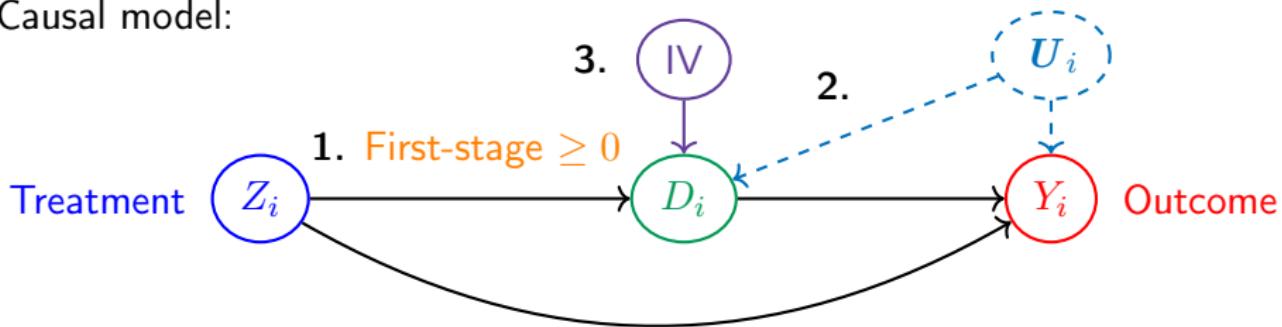
Theorem: Mediation second-stage effects, $Z_i, D_i \rightarrow Y_i$, are identified by the MTE associated Control Functions (CFs).

2. CM with Selection — Identification

MTE assumptions:

- ① Mediator monotonicity
- ② Selection on mediator benefits
- ③ IV for mediator take-up cost.

Causal model:



Proposition: Under MTE assumptions, the mediator MTE is identified.

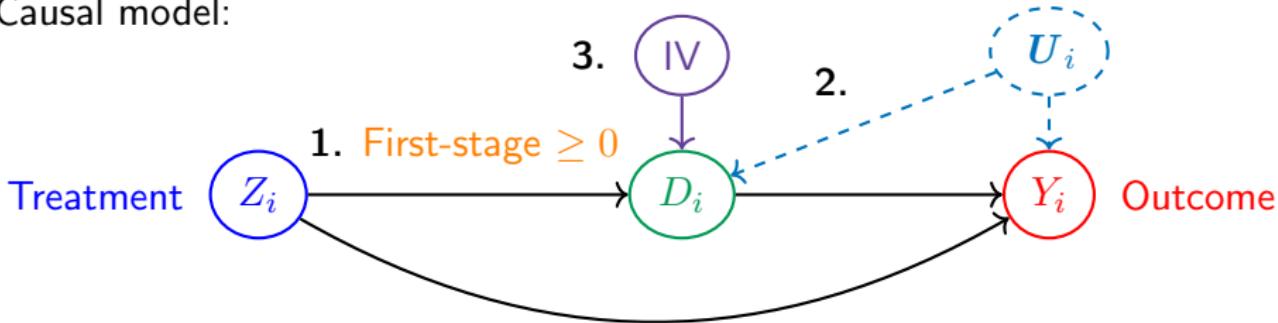
Intuition: Identifies ADE + AIE by extrapolating from IV compliers to mediator compliers (MTE extrapolation e.g., Mogstad Torgovitsky 2018).

2. CM with Selection — Identification

MTE assumptions:

- ① Mediator monotonicity
- ② Selection on mediator benefits
- ③ IV for mediator take-up cost.

Causal model:



Proposition: Under MTE assumptions, the mediator MTE is identified.

Intuition: Identifies ADE + AIE by extrapolating from IV compliers to mediator compliers (MTE extrapolation e.g., Mogstad Torgovitsky 2018).

2. CM with Selection — Estimation

In practice, this means two-stage CM estimation, with CF in second-stage.

Parametric CF Estimation Recipe:

- ① Estimate mediation first-stage with probit, including the IV.
- ② Estimate mediation second-stage by OLS, with Mills ratio CF terms (Heckman 1979).
- ③ Compose CM estimates from two-stage plug-in estimates (same as parametric MTEs, Björklund Moffitt 1987).

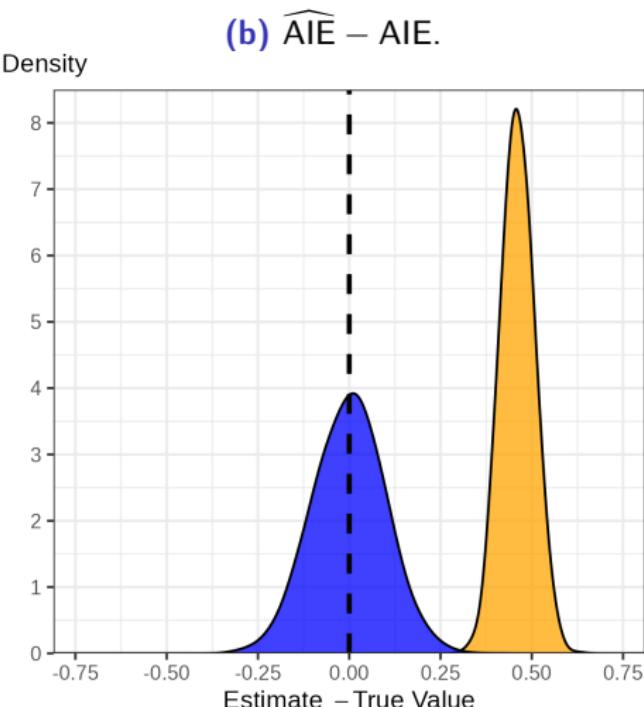
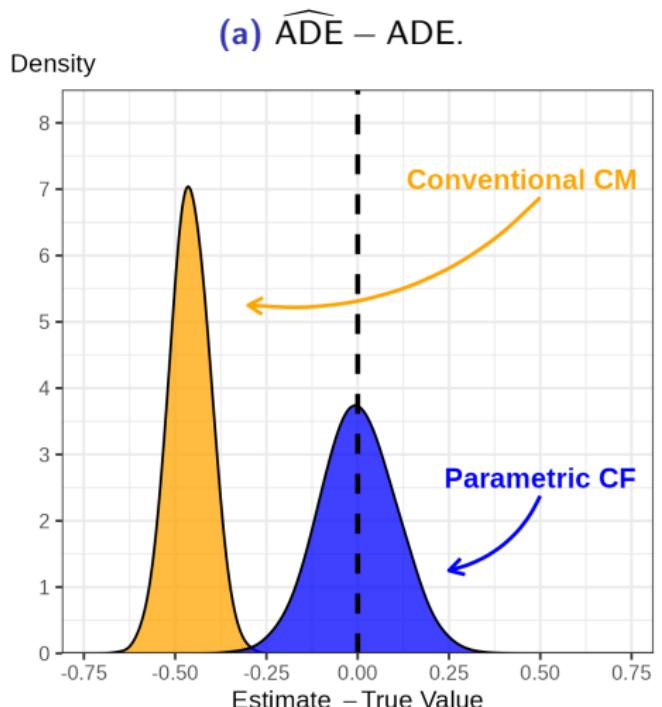
Semi-parametric CF Estimation Recipe:

Replace 2. with semi-parametric CFs (same estimation as MTEs).

⇒ Conventional CM estimates (two-stages) + IV-guided CF adjustment.

2. CM with Selection — Estimation

Figure: CM Estimates from 10,000 DGPs with **Normal** Errors.



2. CM with Selection — Estimation

In practice, this means two-stage CM estimation, with CF in second-stage.

Parametric CF Estimation Recipe:

- ① Estimate mediation first-stage with probit, including the IV.
- ② Estimate mediation second-stage by OLS, with Mills ratio CF terms (Heckman 1979).
- ③ Compose CM estimates from two-stage plug-in estimates (same as parametric MTEs, Björklund Moffitt 1987).

Semi-parametric CF Estimation Recipe:

Replace 2. with semi-parametric CFs (same estimation as MTEs).

⇒ Conventional CM estimates (two-stages) + IV-guided CF adjustment.

2. CM with Selection — Estimation

In practice, this means two-stage CM estimation, with CF in second-stage.

Parametric CF Estimation Recipe:

- ① Estimate mediation first-stage with probit, including the IV.
- ② Estimate mediation second-stage by OLS, with Mills ratio CF terms (Heckman 1979).
- ③ Compose CM estimates from two-stage plug-in estimates (same as parametric MTEs, Björklund Moffitt 1987).

Semi-parametric CF Estimation Recipe:

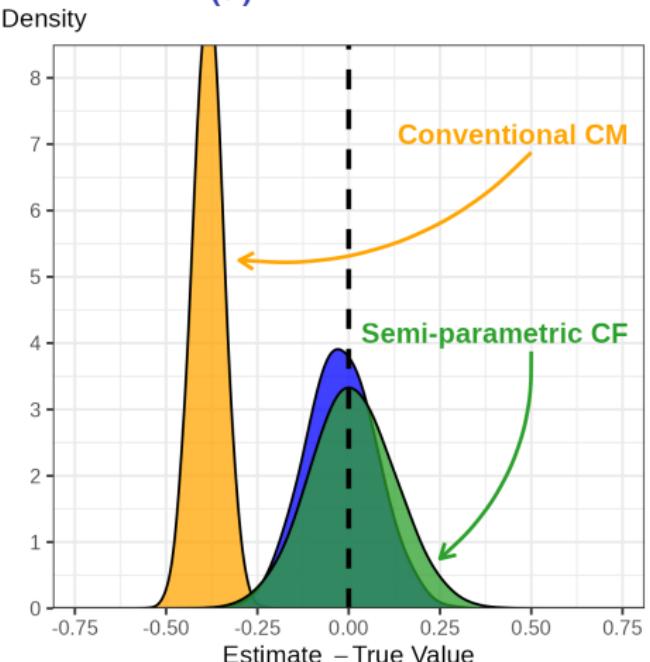
Replace 2. with semi-parametric CFs (same estimation as MTEs).

⇒ Conventional CM estimates (two-stages) + IV-guided CF adjustment.

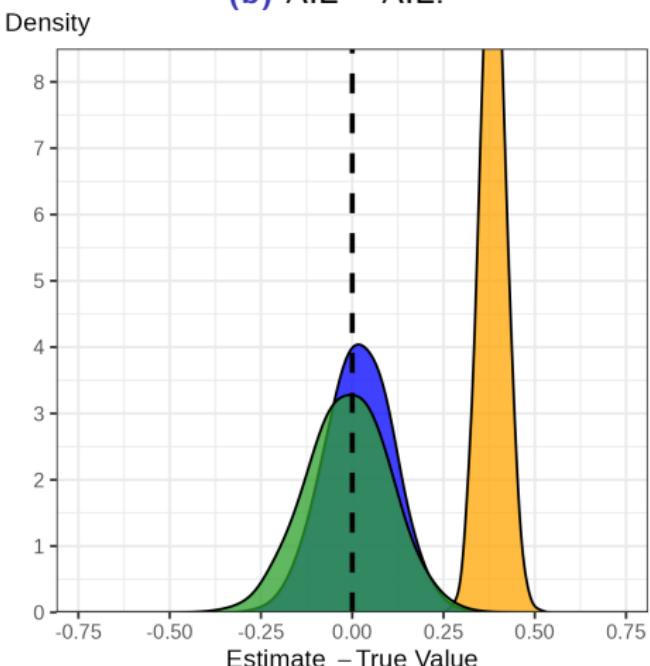
2. CM with Selection — Estimation

Figure: CM Estimates from 10,000 DGPs with **Uniform** Errors.

(a) $\widehat{ADE} - ADE$.

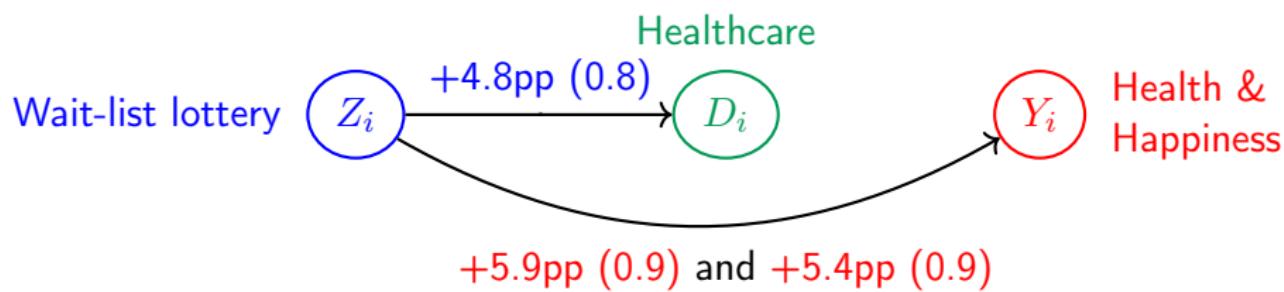


(b) $\widehat{AIE} - AIE$.



3. Returning to Oregon

Winning access to Medicaid increases healthcare usage, and self-reported health & well-being:

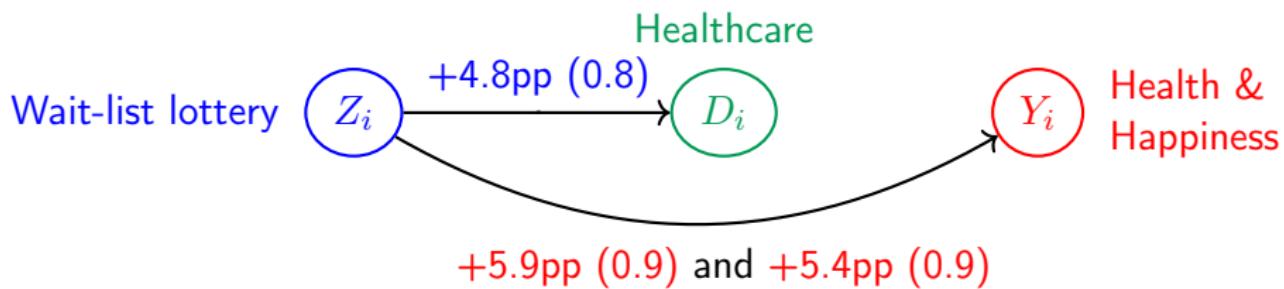


CM is quantitatively estimating the entire system:

- Use correlational estimate of $D_i \rightarrow Y_i$
- Does visiting healthcare at least once increase self-reported health & happiness 12 months later?
- OLS is -1.9pp (1.0) and +3.3 (1.1) for health & happiness, respectively.

3. Returning to Oregon

Winning access to Medicaid increases healthcare usage, and self-reported health & well-being:



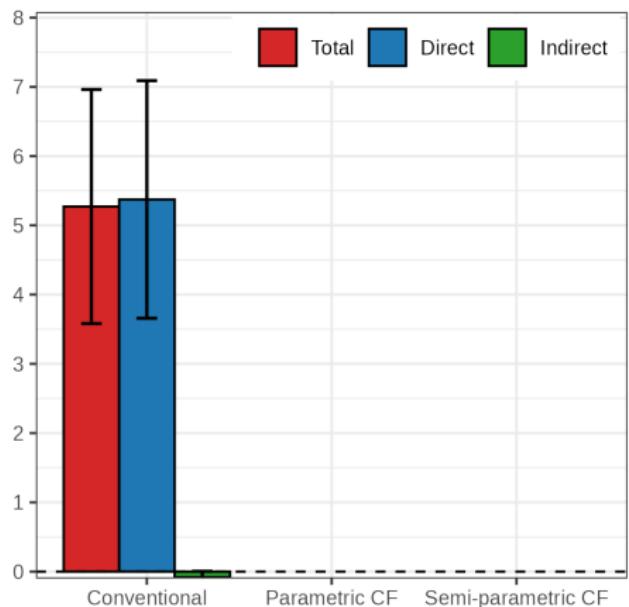
CM is quantitatively estimating the entire system:

- Use correlational estimate of $D_i \rightarrow Y_i$
- Does visiting healthcare at least once increase self-reported health & happiness 12 months later?
- OLS is -1.9pp (1.0) and +3.3 (1.1) for health & happiness, respectively.

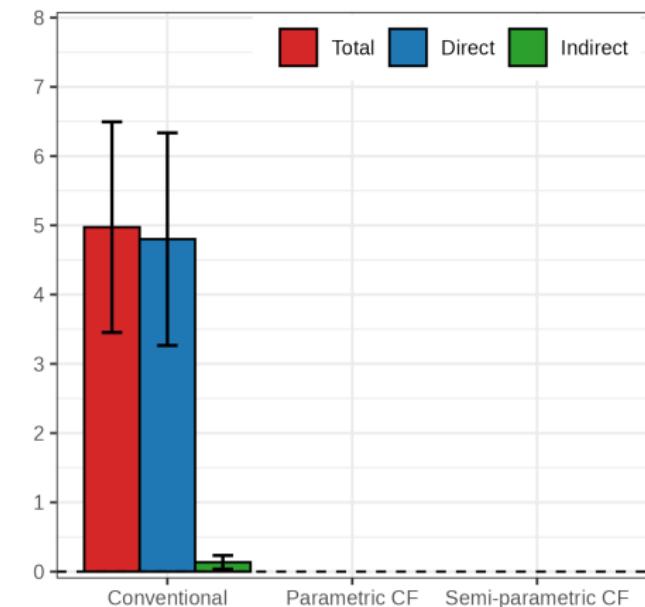
3. Returning to Oregon

Conventional CM estimates lottery effects as mostly direct, ≈ 0 healthcare.

Estimate, percent effect on self-reported health

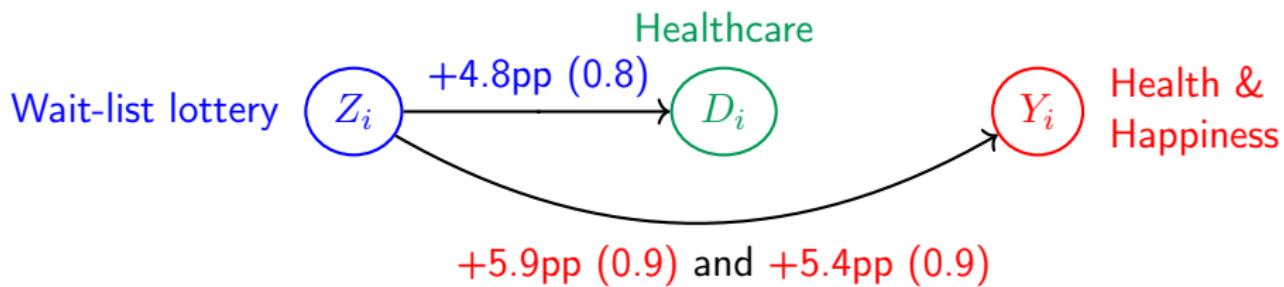


Estimate, percent effect on self-reported happiness



3. Returning to Oregon

Winning access to Medicaid increases healthcare usage, and self-reported health & well-being:

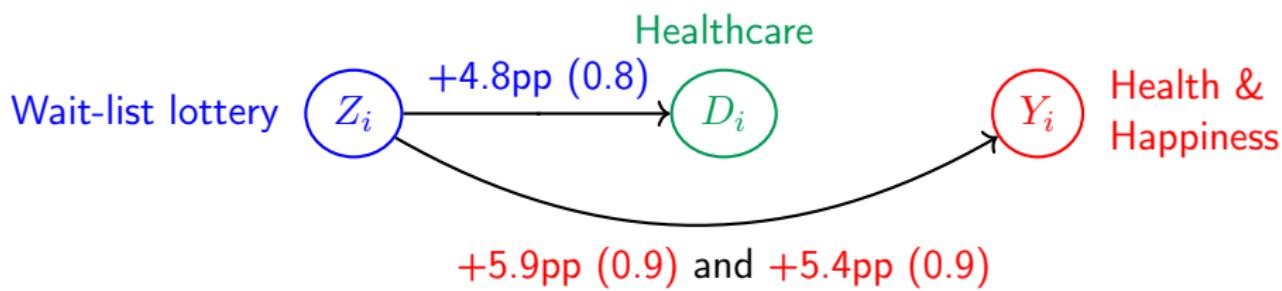


My approach to CM is modelling selection-into- D_i via mediator MTE:

- Uses an estimate of $D_i \rightarrow Y_i$ (plus complier extrapolation)
- Regular healthcare location pre-lottery serves as an excluded IV IV.
- IV estimate of $D_i \rightarrow Y_i$ is +34.2pp (4.4) and +47.2pp (4.5) for health & happiness, respectively.

3. Returning to Oregon

Winning access to Medicaid increases healthcare usage, and self-reported health & well-being:



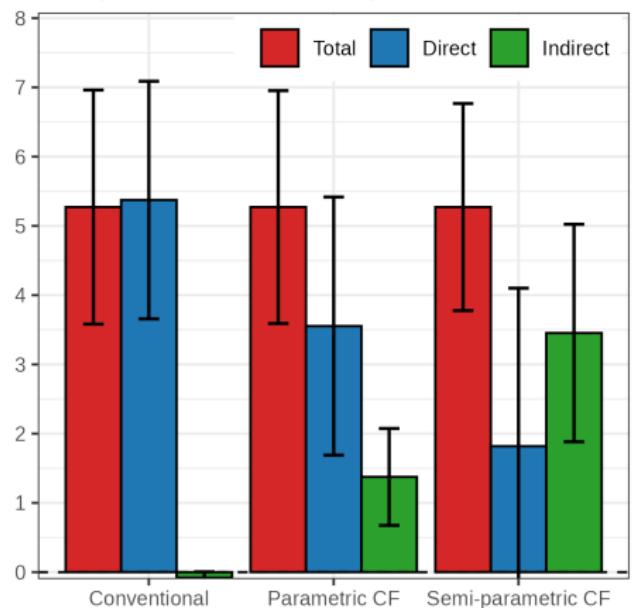
My approach to CM is modelling selection-into- D_i via mediator MTE:

- Uses an estimate of $D_i \rightarrow Y_i$ (plus complier extrapolation)
- Regular healthcare location pre-lottery serves as an excluded IV ▶ IV.
- IV estimate of $D_i \rightarrow Y_i$ is $+34.2\text{pp (4.4)}$ and $+47.2\text{pp (4.5)}$ for health & happiness, respectively.

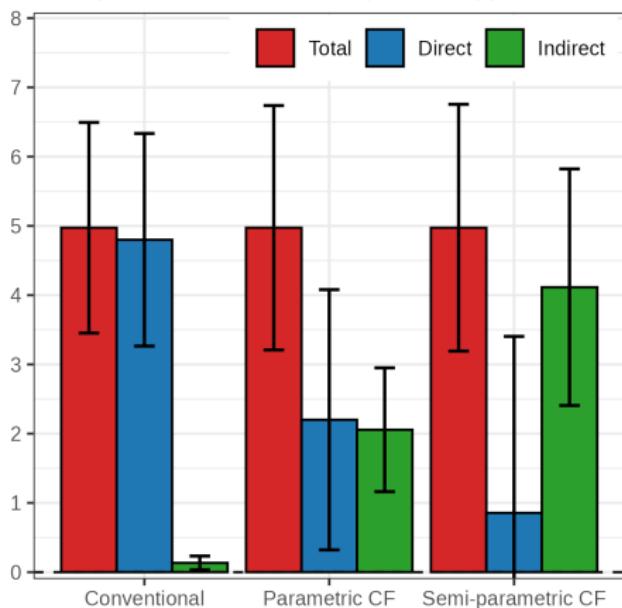
3. Returning to Oregon

Using my approach, with regular healthcare location as an excluded IV, restores indirect effect through increasing healthcare visitation.

Estimate, percent effect on self-reported health



Estimate, percent effect on self-reported happiness



Conclusion

Overview:

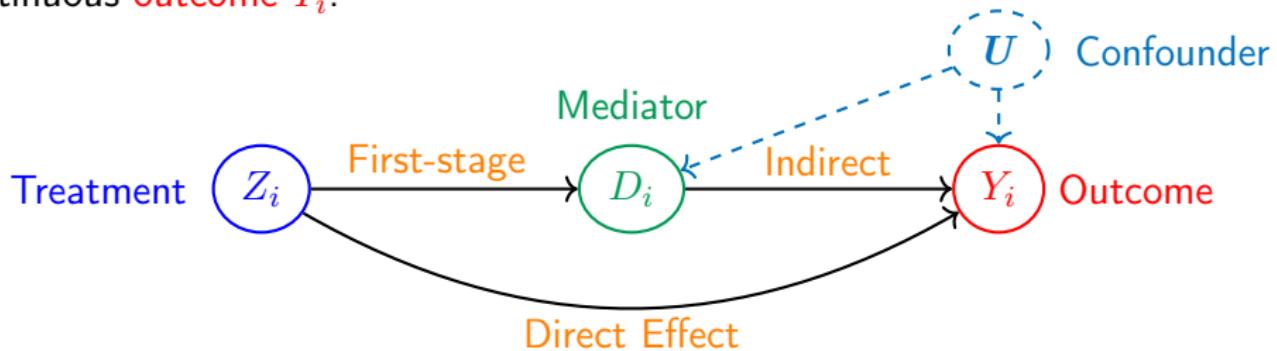
- ① CM as alternative to “suggestive evidence for mechanisms.”
- ② Selection bias in conventional CM analyses with no case for mediator ignorability
 - Noted problems in the most popular methods for CM, pertinent for economic applications.
- ③ Connect CM with labour theory + selection-into-treatment + MTEs
 - Valid CM identification where structural model appropriate.

Caveats and points to remember:

- Structural assumptions and IV for identification + estimation (not ideal).
- Application to Oregon Health Insurance Experiment in the paper, showing health + well-being effects mediated by healthcare (wide confidence intervals).
- **Credible CM analyses are hard in practice.**

Appendix: CM Guiding Model

Consider binary treatment $Z_i = 0, 1$, binary mediator $D_i = 0, 1$, and continuous outcome Y_i .



Average Direct Effect (ADE) : $\mathbb{E} \left[Y_i \left(1, D_i(Z_i) \right) - Y_i \left(0, D_i(Z_i) \right) \right]$

- ADE is causal effect $Z \rightarrow Y$, blocking the indirect D_i path.

Average Indirect Effect (AIE) : $\mathbb{E} \left[Y_i \left(Z_i, D_i(1) \right) - Y_i \left(Z_i, D_i(0) \right) \right]$

- AIE is causal effect of $D_i(Z_i) \rightarrow Y_i$, blocking the direct Z_i path.

Group Difference — ADE

CM effects contaminated by (less interpretable) bias terms.

$$\text{CM Estimand} = \text{ADEM} + \text{Selection Bias}$$

$$\underbrace{\mathbb{E}_{D_i} \left[\mathbb{E}[Y_i | Z_i = 1, D_i] - \mathbb{E}[Y_i | Z_i = 0, D_i] \right]}_{\text{Estimand, Direct Effect}} \\ = \underbrace{\mathbb{E}_{D_i=d'} \left[\mathbb{E}[Y_i(1, D_i(Z_i)) - Y_i(0, D_i(Z_i)) | D_i(1) = d'] \right]}_{\text{Average Direct Effect on Mediator (ADEM) take-up — i.e., } D_i(1) \text{ weighted}} \\ + \underbrace{\mathbb{E}_{D_i} \left[\mathbb{E}[Y_i(0, D_i(Z_i)) | D_i(1) = d'] - \mathbb{E}[Y_i(0, D_i(Z_i)) | D_i(0) = d'] \right]}_{\text{Selection Bias}}$$

The weighted ADE you get here is a positive weighted sum of local ADEs, but with policy irrelevant weights $D_i(1) = d'$.

⇒ consider this group bias, noting difference from true ADE.

[Back](#)

Selection Bias — Direct Effect

CM Effects + contaminating bias.

$$\text{CM Estimand} = \text{ADE} + (\text{Selection Bias} + \text{Group difference bias})$$

► Model

$$\begin{aligned} & \underbrace{\mathbb{E}_{D_i=d'} \left[\mathbb{E} [Y_i | Z_i = 1, D_i = d'] - \mathbb{E} [Y_i | Z_i = 0, D_i = d'] \right]}_{\text{Estimand, Direct Effect}} \\ &= \underbrace{\mathbb{E} [Y_i(1, D_i(Z_i)) - Y_i(0, D_i(Z_i))]}_{\text{Average Direct Effect}} \\ &+ \underbrace{\mathbb{E}_{D_i=d'} \left[\mathbb{E} [Y_i(0, D_i(Z_i)) | D_i(1) = d'] - \mathbb{E} [Y_i(0, D_i(Z_i)) | D_i(0) = d'] \right]}_{\text{Selection Bias}} \\ &+ \underbrace{\mathbb{E}_{D_i=d'} \left[\begin{array}{l} \left(1 - \Pr(D_i(1) = d')\right) \\ \times \left(\mathbb{E} [Y_i(1, D_i(Z_i)) - Y_i(0, D_i(Z_i)) | D_i(1) = 1 - d'] \right. \right. \\ \left. \left. - \mathbb{E} [Y_i(1, D_i(Z_i)) - Y_i(0, D_i(Z_i)) | D_i(0) = d'] \right) \end{array} \right]}_{\text{Group difference bias}} \end{aligned}$$

► Group-diff

Group Difference — AIE

CM effects contaminated by (less interpretable) bias terms.

$$\text{CM Estimand} = \text{AIEM} + (\text{Selection Bias} + \text{Group difference bias})$$

$$\underbrace{\mathbb{E}_{Z_i} \left[\left(\mathbb{E}[D_i | Z_i = 1] - \mathbb{E}[D_i | Z_i = 0] \right) \times \left(\mathbb{E}[Y_i | Z_i, D_i = 1] - \mathbb{E}[Y_i | Z_i, D_i = 0] \right) \right]}_{\text{Estimand, Indirect Effect}}$$

$$= \underbrace{\mathbb{E}[Y_i(Z_i, D_i(1)) - Y_i(Z_i, D_i(0)) | D_i = 1]}_{\text{Average Indirect Effect on Mediated (AIEM) — i.e., } D_i = 1 \text{ weighted}}$$

$$+ \bar{\pi} \underbrace{\left(\mathbb{E}[Y_i(Z_i, 0) | D_i = 1] - \mathbb{E}[Y_i(Z_i, 0) | D_i = 0] \right)}_{\text{Selection Bias}}$$

$$+ \bar{\pi} \underbrace{\left(\frac{1 - \Pr(D_i(1) = 1, D_i(0) = 0)}{\Pr(D_i(1) = 1, D_i(0) = 0)} \right) \left(\begin{aligned} & \mathbb{E}[Y_i(Z_i, 1) - Y_i(Z_i, 0) | D_i(1) = 0 \text{ or } D_i(0)] \\ & - \mathbb{E}[Y_i(Z_i, 1) - Y_i(Z_i, 0)] \end{aligned} \right)}_{\text{Groups difference Bias}}$$

The weighted AIE you get here is not a positive weighted sum of local AI Es, because the AIE is only about $D(Z)$ compliers. [► Model](#).

→ consider this group bias noting difference from true AIE.

[► Back](#)

Selection Bias — Indirect Effect

CM Effects + contaminating bias, where $\bar{\pi} = \Pr(D_i(0) \neq D_i(1))$.

$\text{CM Estimand} = \text{AIE} + \left(\text{Selection Bias} + \text{Group difference bias} \right)$

$$\begin{aligned}
 & \underbrace{\mathbb{E}_{Z_i} \left[\left(\mathbb{E}[D_i | Z_i = 1] - \mathbb{E}[D_i | Z_i = 0] \right) \times \left(\mathbb{E}[Y_i | Z_i, D_i = 1] - \mathbb{E}[Y_i | Z_i, D_i = 0] \right) \right]}_{\text{Estimand, Indirect Effect}} \\
 &= \underbrace{\mathbb{E}[Y_i(Z_i, D_i(1)) - Y_i(Z_i, D_i(0))]}_{\text{Average Indirect Effect}} \\
 &\quad + \bar{\pi} \underbrace{\left(\mathbb{E}[Y_i(Z_i, 0) | D_i = 1] - \mathbb{E}[Y_i(Z_i, 0) | D_i = 0] \right)}_{\text{Selection Bias}} \\
 &\quad + \bar{\pi} \left[\left(1 - \Pr(D_i = 1)\right) \left(\begin{array}{c} \mathbb{E}[Y_i(Z_i, 1) - Y_i(Z_i, 0) | D_i = 1] \\ - \mathbb{E}[Y_i(Z_i, 1) - Y_i(Z_i, 0) | D_i = 0] \end{array} \right) \right. \\
 &\quad \left. + \left(\frac{1 - \Pr(D_i(1) = 1, D_i(0) = 0)}{\Pr(D_i(1) = 1, D_i(0) = 0)} \right) \left(\begin{array}{c} \mathbb{E}[Y_i(Z_i, 1) - Y_i(Z_i, 0) | D_i(Z_i) \neq Z_i] \\ - \mathbb{E}[Y_i(Z_i, 1) - Y_i(Z_i, 0)] \end{array} \right) \right]
 \end{aligned}$$

Groups difference Bias ▶ Group-diff

Semi-parametric Control Functions

Semi-parametric specifications for the CFs λ_0, λ_1 bring some complications to estimating the AIE.

$$\mathbb{E} [Y_i | Z_i, D_i = 0, \mathbf{X}_i] = \alpha + \gamma Z_i + \varphi(\mathbf{X}_i) + \rho_0 \lambda_0(\pi(Z_i; \mathbf{X}_i)),$$

$$\mathbb{E} [Y_i | Z_i, D_i = 1, \mathbf{X}_i] = (\alpha + \beta) + (\gamma + \delta) Z_i + \varphi(\mathbf{X}_i) + \rho_1 \lambda_1(\pi(Z_i; \mathbf{X}_i))$$

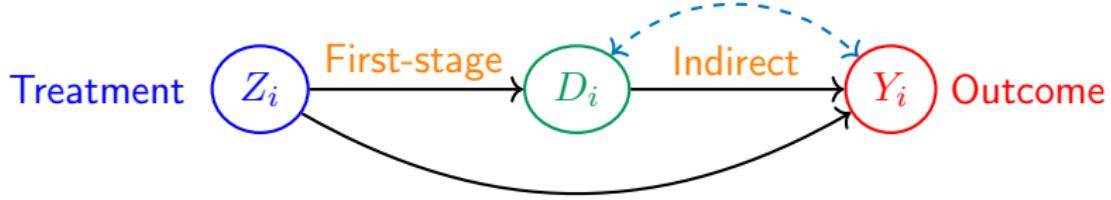
Intercepts, $\alpha, (\alpha + \beta)$, and relevance parameters ρ_0, ρ_1 are not separately identified from the CFs $\lambda_0(.), \lambda_1(.)$ so CF extrapolation term $(\rho_1 - \rho_0)\Gamma(\pi(0; \mathbf{X}_i), \pi(1; \mathbf{X}_i))$ is not directly identified or estimable.

These problems can be avoided by estimating the AIE using its relation to the ATE, $\widehat{\text{AIE}}^{\text{CF}} =$

$$\widehat{\text{ATE}} - (1 - \bar{Z}) \underbrace{\left(\frac{1}{N} \sum_{i=1}^N \widehat{\gamma} + \widehat{\delta} \widehat{\pi}(1; \mathbf{X}_i) \right)}_{\widehat{\text{ADE}} \text{ given } Z_i=1} - \bar{Z} \underbrace{\left(\frac{1}{N} \sum_{i=1}^N \widehat{\gamma} + \widehat{\delta} \widehat{\pi}(0; \mathbf{X}_i) \right)}_{\widehat{\text{ADE}} \text{ given } Z_i=0}.$$

Appendix: CM with Selection

Suppose Z_i is ignorable, D_i is not, so we have the following causal model.



Then this system has the following random coefficient equations:

$$D_i = \phi + \bar{\pi}Z_i + \varphi(\mathbf{X}_i) + U_i$$

$$Y_i = \alpha + \beta D_i + \gamma Z_i + \delta Z_i D_i + \zeta(\mathbf{X}_i) + \underbrace{(1 - D_i)U_{0,i} + D_i U_{1,i}}_{\text{Correlated error term}}$$

where β, γ, δ are functions of $\mu_{d'}(z'; \mathbf{X}_i)$.

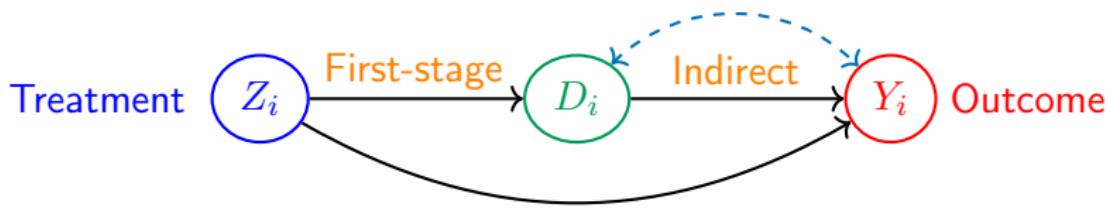
Correlated error term

$$\text{ADE} = \mathbb{E} [\gamma + \delta D_i], \quad \text{AIE} = \mathbb{E} \left[\bar{\pi}(\beta + \delta Z_i + \tilde{U}_i) \right]$$

with $\tilde{U}_i = \mathbb{E} [U_{1,i} - U_{0,i} | \mathbf{X}_i, D_i(0) \neq D_i(1)]$ unobserved complier gains.

Appendix: CM with Selection

Suppose Z_i is ignorable, D_i is not, so we have the following causal model.



Main problem, second-stage is not identified:

$$\begin{aligned}\mathbb{E}[Y_i | Z_i, D_i, \mathbf{X}_i] &= \alpha + \beta D_i + \gamma Z_i + \delta Z_i D_i + \varphi(\mathbf{X}_i) \\ &+ (1 - D_i) \mathbb{E}[U_{0,i} | D_i = 0, \mathbf{X}_i] \\ &+ D_i \mathbb{E}[U_{1,i} | D_i = 1, \mathbf{X}_i]\end{aligned}$$

Unobserved D_i confounding.

Identification intuition: Identify second-stage via MTE control function.

Appendix: CM with Selection — Identification

Assume:

- ① Mediator monotonicity, $\Pr(D_i(0) \leq D_i(1) | \mathbf{X}_i) = 1$
 $\implies D_i(z') = \mathbb{1}\{U_i \leq \pi(z'; \mathbf{X}_i)\}, \text{ for } z' = 0, 1$ (Vycatil 2002).
- ② Selection on mediator benefits, $\text{Cov}(U_i, U_{0,i}), \text{Cov}(U_i, U_{1,i}) \neq 0$
 \implies First-stage take-up informs second-stage confounding.
- ③ There is an IV for the mediator, \mathbf{X}_i^{IV} among control variables \mathbf{X}_i .
 $\implies \pi(Z_i; \mathbf{X}_i) = \Pr(D_i = 1 | Z_i, \mathbf{X}_i)$ is separately identified.

Proposition:

$$\begin{aligned} &\mathbb{E}[Y_i(z', 1) - Y_i(z', 0) | Z_i = z', \mathbf{X}_i, U_i = p'] \\ &= \beta + \delta z' + \mathbb{E}[U_{1,i} - U_{0,i} | \mathbf{X}_i, U_i = p'], \quad \text{for } p' \in (0, 1). \end{aligned}$$

Appendix: CM with Selection — Identification

The marginal effect has corresponding Control Functions (CFs), describing unobserved selection-into- D_i ,

$$\rho_0 \lambda_0(p') = \mathbb{E} [U_{0,i} \mid p' \leq U_i], \quad \rho_1 \lambda_1(p') = \mathbb{E} [U_{1,i} \mid U_i \leq p'].$$

These CFs restore second-stage identification, by extrapolating from \mathbf{X}_i^{IV} compliers to $D_i(Z_i)$ mediator compliers,

$$\begin{aligned} \mathbb{E} [Y_i \mid Z_i, D_i, \mathbf{X}_i] &= \alpha + \beta D_i + \gamma Z_i + \delta Z_i D_i + \varphi(\mathbf{X}_i) \\ &\quad + \underbrace{\rho_0 (1 - D_i) \lambda_0(\pi(Z_i; \mathbf{X}_i)) + \rho_1 D_i \lambda_1(\pi(Z_i; \mathbf{X}_i))}_{\text{CF adjustment.}} \end{aligned}$$

This adjusted second-stage re-identifies the ADE and AIE,

$$\text{ADE} = \mathbb{E} [\gamma + \delta D_i], \quad \text{AIE} = \mathbb{E} \left[\bar{\pi} \left(\beta + \delta Z_i + \underbrace{(\rho_1 - \rho_0) \Gamma(\pi(0; \mathbf{X}_i), \pi(1; \mathbf{X}_i))}_{\text{Mediator compliers extrapolation.}} \right) \right]$$

Appendix: CM with Selection — Estimation

Will explain how estimation works, with simulation evidence.

- ① Random treatment $Z_i \sim \text{Binom}(0.5)$, for $n = 5,000$.
- ② $(U_{0,i}, U_{1,i}) \sim \text{BivariateNormal}(0, 0, \sigma_0, \sigma_1, \rho)$, Costs $C_i \sim N(0, 0.5)$.

Roy selection-into- D_i , with constant partial effects + interaction term.

$$D_i(z') = \mathbb{1} \left\{ C_i \leq Y_i(z', 1) - Y_i(z', 0) \right\},$$
$$Y_i(z', d') = (z' + d' + z'd') + U_{d'} \quad \text{for } z', d' = 0, 1.$$

Following the previous, these data have the following first and second-stage equations, where X_i^{IV} is an additive cost IV:

$$D_i = \mathbb{1} \left\{ C_i - \left(U_{1,i} - U_{0,i} \right) \leq Z_i - X_i^{\text{IV}} \right\}$$
$$Y_i = Z_i + D_i + Z_i D_i + (1 - D_i) U_{0,i} + D_i U_{1,i}.$$

\implies unobserved confounding by BivariateNormal $(U_{0,i}, U_{1,i})$.

Appendix: CM with Selection — Estimation

Errors are normal, so system is Heckman (1979) selection model.

CFs are the inverse Mills ratio, with $\phi(\cdot)$ normal pdf and $\Phi(\cdot)$ normal cdf,

$$\lambda_0(p') = \frac{\phi(-\Phi^{-1}(p'))}{\Phi(-\Phi^{-1}(p'))}, \quad \lambda_1(p') = \frac{\phi(\Phi^{-1}(p'))}{\Phi(\Phi^{-1}(p'))}, \quad \text{for } p' \in (0, 1).$$

Parametric Estimation Recipe:

- ① Estimate first-stage $\pi(Z_i; \mathbf{X}_i)$ with probit, including \mathbf{X}_i^{IV} .
 - ② Include λ_0, λ_1 CFs in second-stage OLS estimation.
 - ③ Compose CM estimates from two-stage plug-in estimates.
-

→ Same as conventional CM estimates (two-stages), with CFs added.

$$\widehat{\text{ADE}} = \mathbb{E} \left[\widehat{\gamma} + \widehat{\delta} D_i \right], \quad \widehat{\text{AIE}} = \mathbb{E} \left[\widehat{\pi} \left(\widehat{\beta} + \widehat{\delta} Z_i + \underbrace{(\widehat{\rho}_1 - \widehat{\rho}_0) \Gamma(\widehat{\pi}(0; \mathbf{X}_i), \widehat{\pi}(1; \mathbf{X}_i))}_{\text{Mediator compliers extrapolation.}} \right) \right]$$

Appendix: CM with Selection — Estimation

If errors are not normal, then CFs do not have a known form, so semi-parametrically estimate them (e.g., splines).

$$\mathbb{E}[Y_i | Z_i, D_i = 0, \mathbf{X}_i] = \alpha + \gamma Z_i + \varphi(\mathbf{X}_i) + \rho_0 \lambda_0(\pi(Z_i; \mathbf{X}_i)),$$

$$\mathbb{E}[Y_i | Z_i, D_i = 1, \mathbf{X}_i] = (\alpha + \beta) + (\gamma + \delta) Z_i + \varphi(\mathbf{X}_i) + \rho_1 \lambda_1(\pi(Z_i; \mathbf{X}_i))$$

Semi-parametric Estimation Recipe:

- ① Estimate first-stage $\pi(Z_i; \mathbf{X}_i)$, including \mathbf{X}_i^{IV} .
- ② Estimate second-stage separately for $D_i = 0$ and $D_i = 1$, with regressors $\lambda_0(p'), \lambda_1(p')$, semi-parametric in $\hat{\pi}(Z_i; \mathbf{X}_i)$.
- ③ Compose CM estimates from two-stage plug-in estimates.

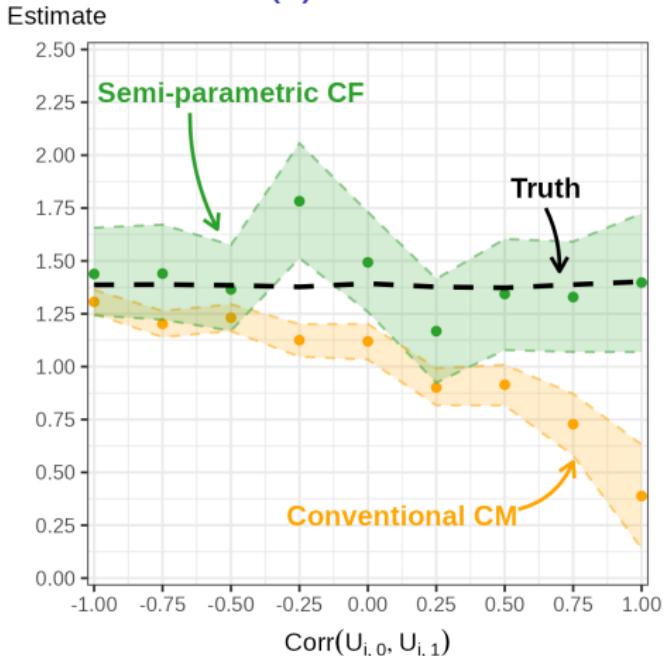
→ Same as conventional CM estimates, with semi-parametric CFs. CFs.

$$\widehat{\text{ADE}} = \mathbb{E} \left[\widehat{\gamma} + \widehat{\delta} D_i \right], \quad \widehat{\text{AIE}} = \mathbb{E} \left[\widehat{\pi} \left(\widehat{\beta} + \widehat{\delta} Z_i + (\widehat{\rho}_1 - \widehat{\rho}_0) \Gamma(\widehat{\pi}(0; \mathbf{X}_i), \widehat{\pi}(1; \mathbf{X}_i)) \right) \right]$$

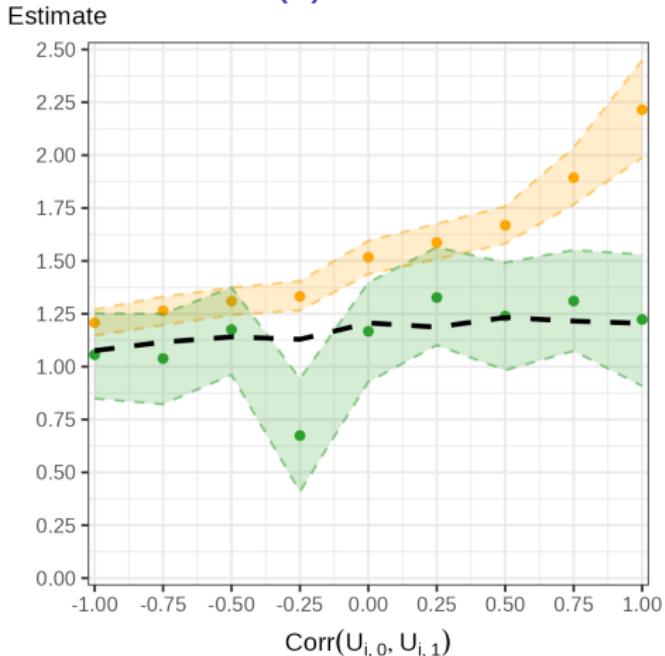
Appendix: CM with Selection — Estimation

Figure: CF Adjusted Estimates Work with Different Error Term Parameters.

(a) ADE.



(b) AIE.



Appendix: OHIE IV

Usual Healthcare Location

