

# Addressing Confounding and Continuous Exposure Measurement Error Using Corrected Score Functions

Brian Richardson

Department of Biostatistics, University of North Carolina at Chapel Hill

August 5, 2025

# Acknowledgements

Bryan Blette, PhD



Vanderbilt University

Peter Gilbert, PhD



University of Washington  
Medical Center

Michael Hudgens, PhD

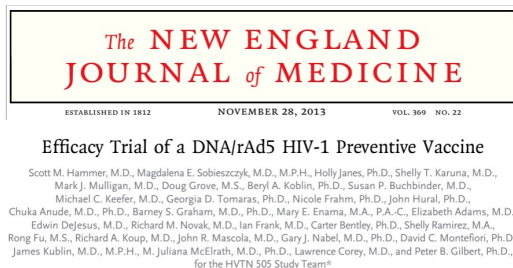


UNC Chapel Hill

This research was supported by the U.S. Public Health Service Grant AI068635, the National Institute of Allergy And Infectious Diseases of the National Institutes of Health (NIH) under Award Number R37 AI054165, and the National Institute of Environmental Health Sciences of the NIH under Award Number T32

# Motivation: HVTN 505 trial

- **HVTN 505 trial:** trial of a preventive HIV vaccine
- stopped early after reaching predetermined cutoffs for efficacy futility [2]



# Motivation: HVTN 505 trial

- several possible biomarker correlates of HIV among vaccine recipients [3, 1, 4]
- of interest to assess the effect of these biomarkers on risk of HIV acquisition
- biomarker-HIV relationship is confounded
- biomarkers are measured with error

*The Journal of Infectious Diseases*

MAJOR ARTICLE



## Higher T-Cell Responses Induced by DNA/rAd5 HIV-1 Preventive Vaccine Are Associated With Lower HIV-1 Infection Risk in an Efficacy Trial

Holly E. James,<sup>1</sup> Kristen W. Cohen,<sup>2</sup> Nicole Frahm,<sup>1</sup> Stephen C. De Rosa,<sup>1</sup> Brittany Sanchez,<sup>1</sup> John Hural,<sup>1</sup> Craig A. Magaret,<sup>1</sup> Shelly Karuna,<sup>1</sup> Carter Bentley,<sup>1</sup> Raphael Gottardo,<sup>1</sup> Greg Finak,<sup>1</sup> Douglas Grove,<sup>1</sup> Mingchao Shen,<sup>1</sup> Barney S. Graham,<sup>1</sup> Richard A. Koup,<sup>3</sup> Mark J. Mulligan,<sup>4</sup> Beryl Koblin,<sup>5</sup> Susan P. Buchbinder,<sup>6</sup> Michael C. Keefer,<sup>7</sup> Elizabeth Adams,<sup>8</sup> Chuka Anude,<sup>1a</sup> Lawrence Corey,<sup>1</sup> Magdalena Sobieszczyk,<sup>1a</sup> Scott M. Hammer,<sup>1a</sup> Peter B. Gilbert,<sup>1</sup> and M. Juliana McElrath<sup>1</sup>

# Goal

To estimate the **causal effect** of a continuous exposure on an outcome when  
(i) the exposure-outcome relationship is potentially **confounded**

# Goal

To estimate the **causal effect** of a continuous exposure on an outcome when

- (i) the exposure-outcome relationship is potentially **confounded**
- (ii) the exposure is **measured with error**

# What do we Mean by Causal Effect?

“What would be the risk of HIV if, possibly counter to fact, somebody were to have biomarker level ***a***?”

# What do we Mean by Causal Effect?

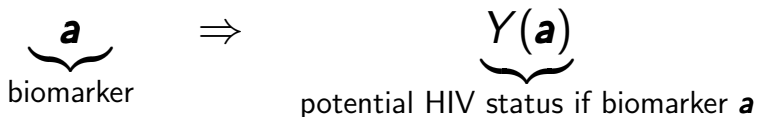
“What would be the risk of HIV if, possibly counter to fact, somebody were to have biomarker level  $a$ ?”

$a$   
biomarker



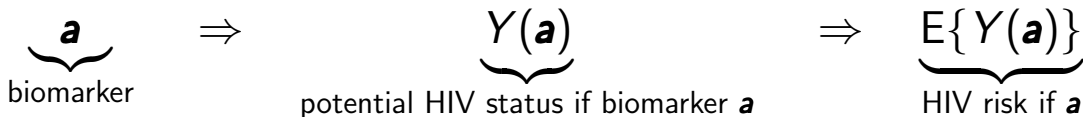
# What do we Mean by Causal Effect?

“What would be the risk of HIV if, possibly counter to fact, somebody were to have biomarker level  $\mathbf{a}$ ?”



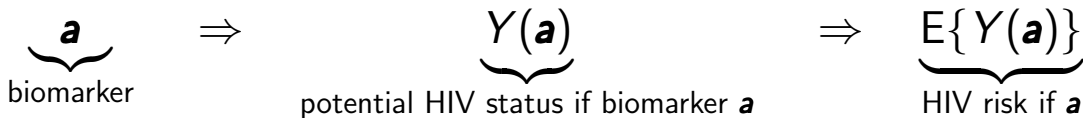
# What do we Mean by Causal Effect?

“What would be the risk of HIV if, possibly counter to fact, somebody were to have biomarker level  $\mathbf{a}$ ?”



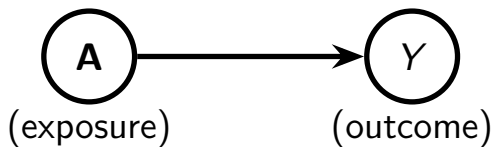
# What do we Mean by Causal Effect?

“What would be the risk of HIV if, possibly counter to fact, somebody were to have biomarker level  $\mathbf{a}$ ?”

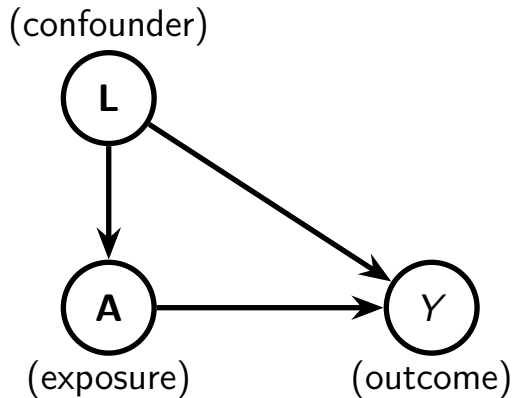


Estimand (**dose-response surface**):  $\eta(\mathbf{a}) \equiv E[Y(\mathbf{a})]$  for  $\mathbf{a} \in \mathcal{A}$

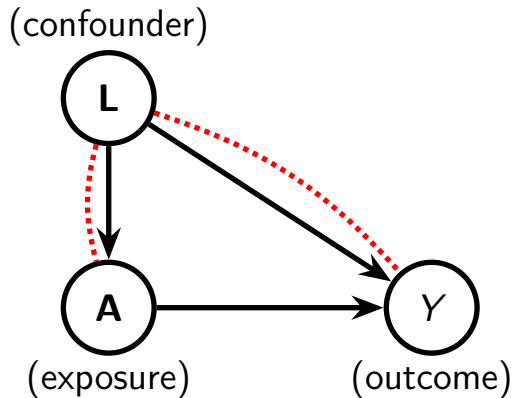
# Confounding



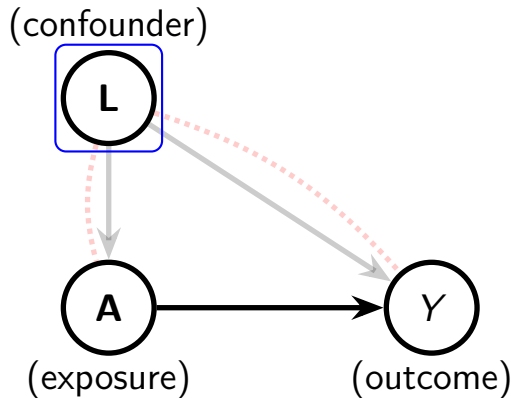
# Confounding



# Confounding



# Confounding



# Addressing Confounding (3 Ways)

- G-Formula (GF)
- Inverse Probability Weighting (IPW)
- Doubly Robust Methods (DR)



# Addressing Confounding (3 Ways)

- G-Formula (GF)
  - ▶ Fit an **outcome regression model** for  $\mu(\mathbf{L}, \mathbf{A}; \beta) \equiv E(Y|\mathbf{L}, \mathbf{A})$
  - ▶ Marginalize over the distribution of confounders:  $\hat{\eta}(\mathbf{a}) = n^{-1} \sum_{i=1}^n \mu(\mathbf{L}_i, \mathbf{a}; \hat{\beta})$
- Inverse Probability Weighting (IPW)
- Doubly Robust Methods (DR)

# Addressing Confounding (3 Ways)

- G-Formula (GF)
- Inverse Probability Weighting (IPW)
  - ▶ Fit a **propensity model** for the distribution of  $\mathbf{A}|\mathbf{L}$
  - ▶ Weight each observation based on its propensity score
  - ▶ Fit a regression model for  $Y$  on  $\mathbf{A}$  using weighted observations
- Doubly Robust Methods (DR)

# Addressing Confounding (3 Ways)

- G-Formula (GF)
- Inverse Probability Weighting (IPW)
- Doubly Robust Methods (DR)
  - ▶ Fit a **propensity model** for  $\mathbf{A}$  given  $\mathbf{L}$
  - ▶ Weight each observation based on its propensity score
  - ▶ Fit an **outcome regression model** for  $\mu(\mathbf{L}, \mathbf{A}; \beta) \equiv E(Y|\mathbf{L}, \mathbf{A})$  using weighted observations
  - ▶ Marginalize over the distribution of confounders:  $\hat{\eta}(\mathbf{a}) = n^{-1} \sum_{i=1}^n \mu(\mathbf{L}_i, \mathbf{a}; \hat{\beta})$

# Estimating Equations

**Estimating Function:** a function of the observed data and the parameter of interest

$$\underbrace{\psi}_{\text{est. fun.}} \left( \quad ; \quad \right)$$

# Estimating Equations

**Estimating Function:** a function of the observed data and the parameter of interest

$$\underbrace{\psi}_{\text{est. fun.}} \left( \underbrace{Y, L, A}_{\text{data}} ; \right)$$

# Estimating Equations

**Estimating Function:** a function of the observed data and the parameter of interest

$$\underbrace{\psi}_{\text{est. fun.}} \left( \underbrace{Y, L, A}_{\text{data}} ; \underbrace{\theta}_{\text{param.}} \right)$$

# Estimating Equations

**Estimating Function:** a function of the observed data and the parameter of interest

$$\underbrace{\psi}_{\text{est. fun.}} \left( \underbrace{Y, L, A}_{\text{data}} ; \underbrace{\theta}_{\text{param.}} \right)$$

that is **unbiased**, meaning

# Estimating Equations

**Estimating Function:** a function of the observed data and the parameter of interest

$$\underbrace{\Psi}_{\text{est. fun.}} \left( \underbrace{Y, L, A}_{\text{data}} ; \underbrace{\theta}_{\text{param.}} \right)$$

that is **unbiased**, meaning

$$E\{\Psi(Y, L, A; \underbrace{\theta_0}_{\text{true}})\} = \mathbf{0}.$$



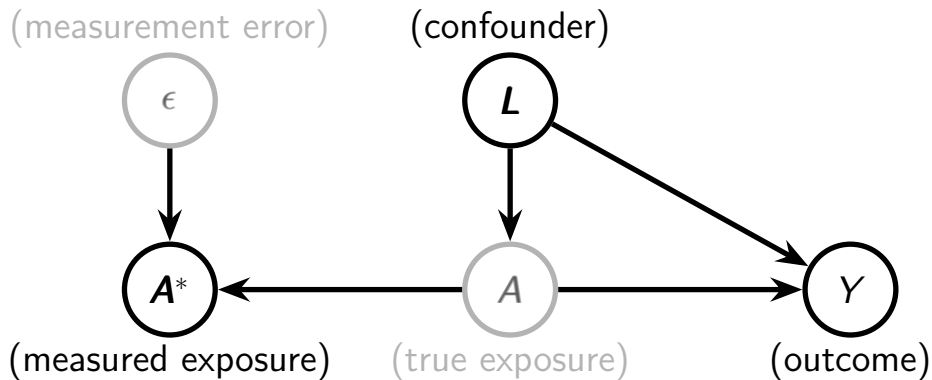
# Estimating Equations (contd)

- Given an estimating function  $\Psi$  and observed data  $\{(Y_i, L_i, A_i) : i = 1, \dots, n\}$ , we can find an estimator  $\hat{\theta}$  as the solution to

$$\sum_{i=1}^n \underbrace{\Psi(Y_i, L_i, A_i; \theta)}_{\text{observed}} = \mathbf{0}.$$

- $\hat{\theta}$  is consistent and asymptotically normal and has a simple variance estimator [6].
- The G-formula, IPW, and DR estimators are all solutions to estimating equations with functions  $\Psi_{GF}, \Psi_{IPW}, \Psi_{DR}$ .

# Measurement Error



# Addressing Confounding and Measurement Error

Can we just substitute  $\mathbf{A}^*$  for  $\mathbf{A}$  and find the solution to

$$\sum_{i=1}^n \Psi(Y_i, L_i, \underbrace{\mathbf{A}_i^*}_{\text{mismeasured}}; \theta) = \mathbf{0}?$$

No! This leads to bias in  $\hat{\theta}$  because

$$E\{\Psi(Y, L, \mathbf{A}^*; \theta_0)\} \neq \mathbf{0}.$$

We need a new estimating function  $\tilde{\Psi}$  such that

$$E\left\{ \underbrace{\tilde{\Psi}}_{\text{new est. fun.}}(Y, L, \underbrace{\mathbf{A}^*}_{\text{mismeasured}}; \theta_0) \right\} = \mathbf{0}.$$

# Corrected Score Functions

Given the “oracle” estimating function  $\Psi$ , the “corrected score” function  $\Psi_{CS}$  can be created following Novick and Stefanski [5]:

- 1 add additional *imaginary* measurement error to the mismeasured exposure:  
 $\tilde{\mathbf{A}} = \mathbf{A}^* + i\tilde{\epsilon}.$
- 2 Plug  $\tilde{\mathbf{A}}$  into  $\Psi$ , and keep only the real part of the complex-valued function.
- 3 Take the expectation over the additional measurement error  $\tilde{\epsilon}$ .

$$\Psi_{CS}(Y, L, \mathbf{A}^*; \theta) = E \left[ \text{Re} \left\{ \Psi_0(Y, L, \tilde{\mathbf{A}}; \theta) \right\} \mid Y, L, \mathbf{A}^* \right]$$

## Corrected Score Functions (contd)

Under certain conditions, the corrected score function  $\Psi_{CS}$  is then unbiased, meaning

$$E\{\Psi_{CS}(Y, L, \underbrace{A^*}_{\text{mismeasured}}; \theta_0)\} = \mathbf{0}.$$

The G-Formula, IPW, and DR estimating functions all satisfy these conditions, and so can be “corrected.”

$$\Psi_{GF} \longrightarrow \Psi_{CS-GF}$$

$$\Psi_{IPW} \longrightarrow \Psi_{CS-IPW}$$

$$\Psi_{DR} \longrightarrow \Psi_{CS-DR}$$

# Monte Carlo Corrected Score Functions

- Sometimes we can find a closed-form algebraic expression for

$$\boldsymbol{\Psi}_{CS}(Y, \mathbf{L}, \mathbf{A}^*; \boldsymbol{\theta}) = \underbrace{\mathbb{E}[\operatorname{Re}\{\boldsymbol{\Psi}_0(Y, \mathbf{L}, \mathbf{A}^* + i\tilde{\boldsymbol{\epsilon}}; \boldsymbol{\theta})\} \mid Y, \mathbf{L}, \mathbf{A}^*]}_{\mathbb{E}\{f(\tilde{\boldsymbol{\epsilon}})\}}.$$

# Monte Carlo Corrected Score Functions

- Sometimes we can find a closed-form algebraic expression for

$$\boldsymbol{\Psi}_{CS}(Y, \mathbf{L}, \mathbf{A}^*; \boldsymbol{\theta}) = \underbrace{\mathbb{E}[\text{Re}\{\boldsymbol{\Psi}_0(Y, \mathbf{L}, \mathbf{A}^* + i\tilde{\boldsymbol{\epsilon}}; \boldsymbol{\theta})\} \mid Y, \mathbf{L}, \mathbf{A}^*]}_{\mathbb{E}\{f(\tilde{\boldsymbol{\epsilon}})\}}.$$

- Alternatively, we can approximate this expectation with Monte Carlo replicates

# Monte Carlo Corrected Score Functions

- Sometimes we can find a closed-form algebraic expression for

$$\boldsymbol{\Psi}_{CS}(Y, \mathbf{L}, \mathbf{A}^*; \boldsymbol{\theta}) = \underbrace{\mathbb{E}[\operatorname{Re}\{\boldsymbol{\Psi}_0(Y, \mathbf{L}, \mathbf{A}^* + i\tilde{\boldsymbol{\epsilon}}; \boldsymbol{\theta})\} \mid Y, \mathbf{L}, \mathbf{A}^*]}_{\mathbb{E}\{f(\tilde{\boldsymbol{\epsilon}})\}}.$$

- Alternatively, we can approximate this expectation with Monte Carlo replicates

$$\mathbb{E}\{f(\tilde{\boldsymbol{\epsilon}})\} \approx B^{-1} \sum_{b=1}^B f(\tilde{\boldsymbol{\epsilon}}_b)$$



# Monte Carlo Corrected Score Functions

- Sometimes we can find a closed-form algebraic expression for

$$\boldsymbol{\Psi}_{CS}(Y, \mathbf{L}, \mathbf{A}^*; \boldsymbol{\theta}) = \underbrace{\mathbb{E}[\operatorname{Re}\{\boldsymbol{\Psi}_0(Y, \mathbf{L}, \mathbf{A}^* + i\tilde{\boldsymbol{\epsilon}}; \boldsymbol{\theta})\} \mid Y, \mathbf{L}, \mathbf{A}^*]}_{\mathbb{E}\{f(\tilde{\boldsymbol{\epsilon}})\}}.$$

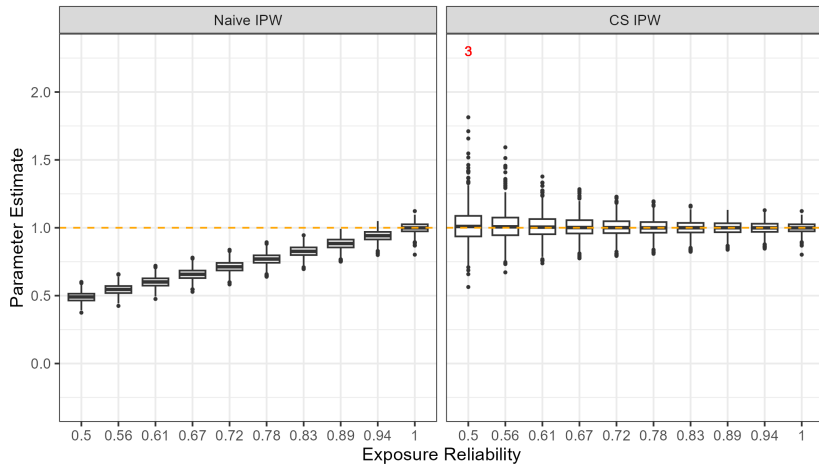
- Alternatively, we can approximate this expectation with Monte Carlo replicates

$$\begin{aligned}\mathbb{E}\{f(\tilde{\boldsymbol{\epsilon}})\} &\approx B^{-1} \sum_{b=1}^B f(\tilde{\boldsymbol{\epsilon}}_b) \\ \implies \boldsymbol{\Psi}_{MCCS}^B(Y, \mathbf{L}, \mathbf{A}^*; \boldsymbol{\theta}) &= B^{-1} \sum_{b=1}^B \operatorname{Re}\{\boldsymbol{\Psi}_0(Y, \mathbf{L}, \mathbf{A}^* + i\tilde{\boldsymbol{\epsilon}}_b; \boldsymbol{\theta})\}\end{aligned}$$

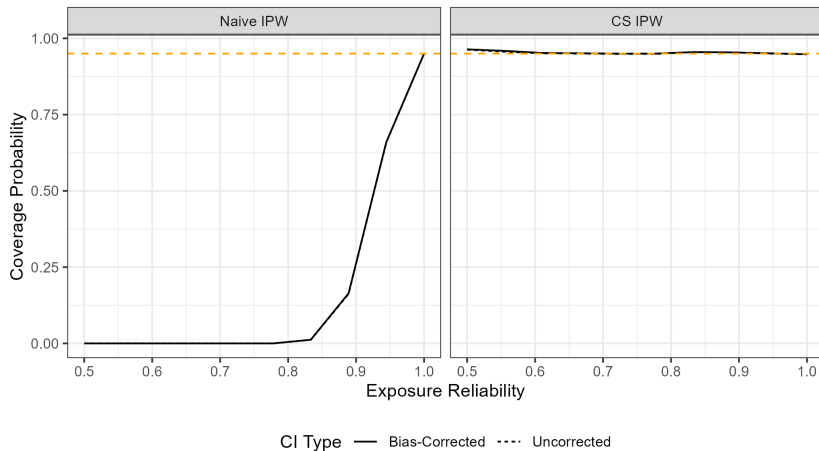
# Simulation Setting

- confounder  $L \sim \mathcal{N}(0, 0.36)$
- exposure  $\mathbf{A} = (A_1, A_2)$  with  $\mathbf{A}|L \sim \mathcal{N}_2(\mathbf{0}, I)$
- exposure measurement error  $\epsilon \sim \mathcal{N}_2(\mathbf{0}, \sigma_{me}^2 I)$
- outcome  $Y$  with  $Y|L, \mathbf{A} \sim \mathcal{N}(A_1 + A_2 + L, 1)$
- implied MSM of  $\eta(\mathbf{a}; \boldsymbol{\gamma}) = \gamma_0 + \gamma_1 a_1 + \gamma_2 a_2$  for  $\boldsymbol{\gamma} = (\gamma_0, \gamma_1, \gamma_2) = (0, 1, 1)$
- sample size  $n = 800$

# Simulation Results: Estimator



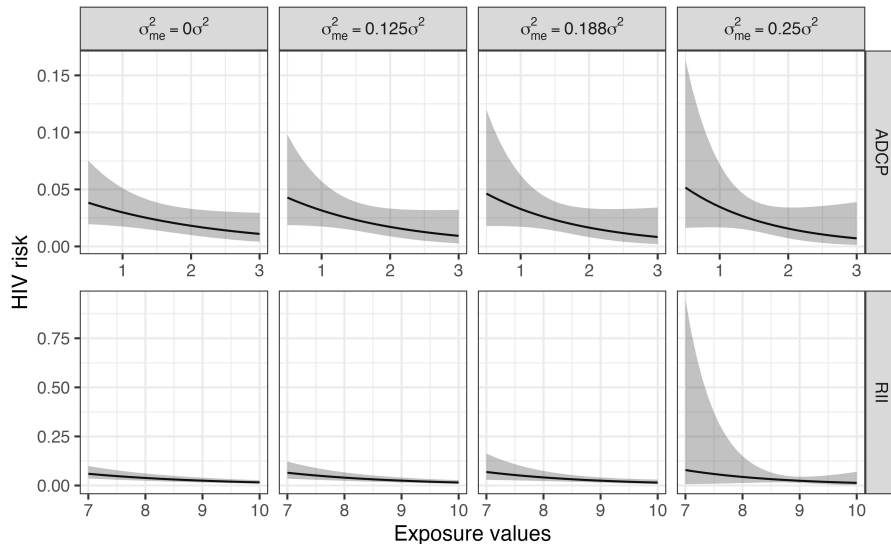
# Simulation Results: Confidence Interval



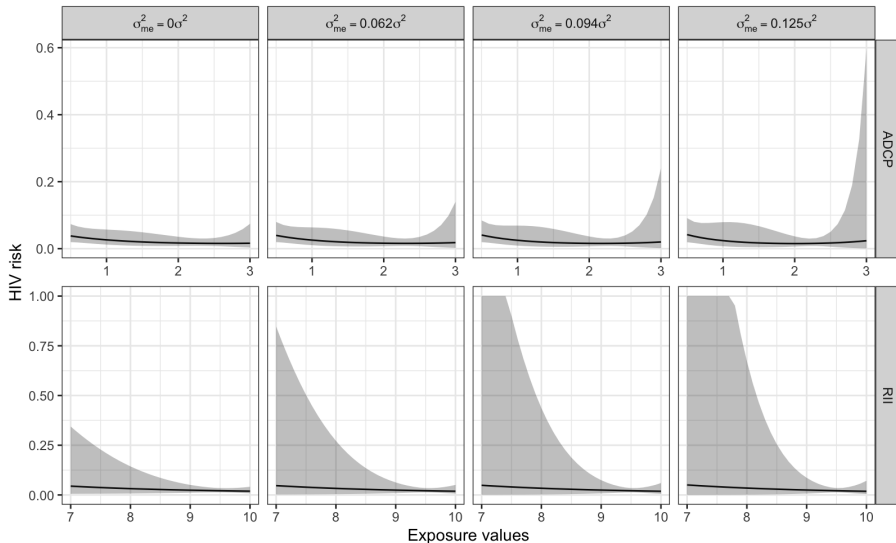
# Application: HVTN 505 Trial

- **two exposures:**
  - (i) antibody-dependent cellular phagocytosis (ADCP)
  - (ii) recruitment of  $Fc\gamma RIIa$  of the H131-Con S gp140 protein (RII)
- **case-cohort sampling:** immunologic markers only measured in stratified random sample of controls
- **covariates:** age, race, BMI, behavior risk, CD4-P, and CD8-P
- **two analyses:**
  - (i) DR estimator with a linear outcome model
  - (ii) g-formula with a quadratic outcome model

# Application: DR Method with Linear Outcome Model



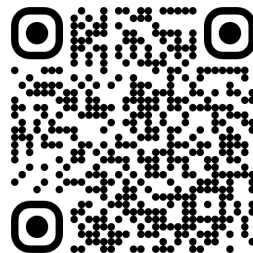
# Application: G-Formula with Quadratic Outcome Model



# Mismex: Causal Inference for Mismeasured Exposures



Paper in *Biometrics*



GitHub R package



# References

- [1] Youyi Fong, Xiaoying Shen, Vicki C Ashley, Aaron Deal, Kelly E Seaton, Chenchu Yu, Shannon P Grant, Guido Ferrari, Allan C deCamp, Robert T Bailer, et al. Modification of the association between T-cell immune responses and human immunodeficiency virus type 1 infection risk by vaccine-induced antibody responses in the HVTN 505 trial. *The Journal of Infectious Diseases*, 217(8):1280–1288, 2018.
- [2] Scott M Hammer, Magdalena E Sobieszczek, Holly Janes, Shelly T Karuna, Mark J Mulligan, Doug Grove, Beryl A Koblin, Susan P Buchbinder, Michael C Keefer, Georgia D Tomaras, et al. Efficacy trial of a DNA/rAd5 HIV-1 preventive vaccine. *New England Journal of Medicine*, 369(22):2083–2092, 2013.
- [3] Holly E Janes, Kristen W Cohen, Nicole Frahm, Stephen C De Rosa, Brittany Sanchez, John Hural, Craig A Magaret, Shelly Karuna, Carter Bentley, Raphael Gottardo, et al. Higher T-cell responses induced by DNA/rAd5 HIV-1 preventive vaccine are associated with lower HIV-1 infection risk in an efficacy trial. *The Journal of Infectious Diseases*, 215(9):1376–1385, 2017.
- [4] Scott D Neidich, Youyi Fong, Shuying S Li, Daniel E Geraghty, Brian D Williamson, William Chad Young, Derrick Goodman, Kelly E Seaton, Xiaoying Shen, Sheetal Sawant, et al. Antibody Fc effector functions and IgG3 associate with decreased HIV-1 risk. *Journal of Clinical Investigation*, 129(11):4838–4849, 2019.
- [5] Steven J Novick and Leonard A Stefanski. Corrected score estimation via complex variable simulation extrapolation. *Journal of the American Statistical Association*, 97(458):472–481, June 2002. ISSN 0162-1459, 1537-274X. doi: 10.1198/016214502760047005. URL <http://www.tandfonline.com/doi/abs/10.1198/016214502760047005>.
- [6] Leonard A Stefanski and Dennis D Boos. The calculus of M-estimation. *The American Statistician*, 56(1):29–38, 2002.

# Appendix

## Appendix: Notation

- true exposure:  $\mathbf{A} = (A_1, \dots, A_m)$
- measured exposure:  $\mathbf{A}^* = (A_1^*, \dots, A_m^*) = \mathbf{A} + \boldsymbol{\epsilon}$
- measurement error:  $\boldsymbol{\epsilon}$
- potential outcome:  $Y(\mathbf{a})$
- observed outcome:  $Y$
- confounders:  $\mathbf{L} = (L_1, L_2, \dots, L_p)$

**Observe:** iid copies of  $(Y_i, \mathbf{L}_i, \mathbf{A}_i^*)$ .

**Estimand:** dose-response curve  $\eta(\mathbf{a}) \equiv E[Y(\mathbf{a})]$  for  $\mathbf{a} \in \mathcal{A}$ .

## Appendix: Assumptions

- (i) **causal consistency**:  $Y = Y(\mathbf{a})$  when  $\mathbf{A} = \mathbf{a}$
- (ii) **conditional exchangeability**:  $Y(\mathbf{a}) \perp\!\!\!\perp \mathbf{A} | \mathbf{L}$  for all  $\mathbf{a} \in \mathcal{A}$
- (iii) **positivity**:  $f_{\mathbf{A}|\mathbf{L}}(\mathbf{a}|\mathbf{l}) > 0$  for all  $\mathbf{l}$  such that  $f_{\mathbf{L}}(\mathbf{l}) > 0$  and for all  $\mathbf{a} \in \mathcal{A}$
- (iv) **independent measurement error**:  $\epsilon \perp\!\!\!\perp (Y, \mathbf{L}, \mathbf{A})$
- (v) **classical additive measurement error**:  $\epsilon \sim \mathcal{N}_m(0, \Sigma_{me})$

## Appendix: G-Formula

- fit the **outcome model**  $\mu(\mathbf{L}, \mathbf{A}; \beta) \equiv E(Y | \mathbf{L}, \mathbf{A})$
- estimate the dose-response curve by marginalizing over the distribution of confounders:  $\hat{\eta}(\mathbf{a}) = n^{-1} \sum_{i=1}^n \mu(\mathbf{L}_i, \mathbf{a}; \hat{\beta})$
- This can be expressed as an M-estimator with estimating function

$$\Psi_{0-GF}(Y, \mathbf{L}, \mathbf{A}; \theta_{GF}) = \begin{bmatrix} \{Y - \mu(\mathbf{L}, \mathbf{A}; \beta)\} \partial_{\beta} \mu(\mathbf{L}, \mathbf{A}; \beta) \\ \eta(\mathbf{a}) - \mu(\mathbf{L}, \mathbf{a}; \beta) \end{bmatrix}$$

## Appendix: IPW

- obtain/estimate **standardized propensity score weights**

$$SW(L, \mathbf{A}) = \frac{f_{\mathbf{A}}(\mathbf{A})}{f_{\mathbf{A}|L}(\mathbf{A}|L)}$$

- use weighted observations to estimate the dose-response curve  $\eta(\mathbf{a}; \gamma)$
- This can be expressed as an M-estimator with estimating function

$$\psi_{0-IPW}(Y, L, \mathbf{A}; \theta_{IPW}) = \begin{bmatrix} \psi_{PS}(L, \mathbf{A}) \\ SW(L, \mathbf{A}) \{Y - \eta(\mathbf{A}; \gamma)\} \partial_{\gamma} \eta(\mathbf{A}; \gamma) \end{bmatrix}$$

## Appendix: DR

- obtain/estimate **standardized propensity score weights**  $SW(\mathbf{L}, \mathbf{A})$
- use weighted observations to estimate the **outcome model**  $\mu(\mathbf{L}, \mathbf{A}; \beta) \equiv E(Y|\mathbf{L}, \mathbf{A})$
- estimate the dose-response curve by marginalizing over the distribution of confounders
- This can be expressed as an M-estimator with estimating function

$$\Psi_{0-DR}(Y, \mathbf{L}, \mathbf{A}; \theta_{DR}) = \begin{bmatrix} \Psi_{PS}(\mathbf{L}, \mathbf{A}) \\ SW(\mathbf{L}, \mathbf{A})\{Y - \mu(\mathbf{L}, \mathbf{A}; \beta)\}\partial_{\beta}\mu(\mathbf{L}, \mathbf{A}; \beta) \\ \eta(\mathbf{a}) - \mu(\mathbf{L}, \mathbf{a}; \beta) \end{bmatrix}$$

- doubly robust\* to models for  $\mu(\mathbf{L}, \mathbf{A}; \beta)$  and  $f_{\mathbf{A}|\mathbf{L}}(\mathbf{A}|\mathbf{L})$ .

## Appendix: Corrected Score Functions

- Suppose the oracle estimating function is **conditionally unbiased**, meaning

$$E\{\Psi_0(Y, L, \mathbf{A}; \theta) | \mathbf{A}\} = 0.$$

- Define the corrected score function as

$$\Psi_{CS}(Y, L, \mathbf{A}^*; \theta) = E \left[ \text{Re} \left\{ \Psi_0(Y, L, \tilde{\mathbf{A}}; \theta) \right\} \mid Y, L, \mathbf{A}^* \right],$$

where  $\tilde{\mathbf{A}} = \mathbf{A}^* + i\tilde{\epsilon}$ ,  $i = \sqrt{-1}$ ,  $\text{Re}(\cdot)$  denotes the real component of a complex number, and  $\tilde{\epsilon} \sim \mathcal{N}(\mathbf{0}, \Sigma_{me})$ .

- Then

$$\begin{aligned} E\{\Psi_{CS}(Y, L, \mathbf{A}^*; \theta) \mid Y, L, \mathbf{A}\} &= \Psi_0(Y, L, \mathbf{A}; \theta) \\ \implies E[E\{\Psi_{CS}(Y, L, \mathbf{A}^*; \theta) \mid Y, L, \mathbf{A}\}] &= E\{\Psi_0(Y, L, \mathbf{A}; \theta)\} \\ \implies E\{\Psi_{CS}(Y, L, \mathbf{A}^*; \theta)\} &= 0 \end{aligned}$$