

Nima S. Hejazi<sup>1</sup> Bryan E. Shepherd <sup>3</sup> Giovanni Parmigiani<sup>1</sup> Keith Barnatchez<sup>1</sup> Kevin P. Josey<sup>2</sup> <sup>2</sup>Colorado School of Public Health <sup>1</sup>Harvard T.H. Chan School of Public Health <sup>3</sup>Vanderbilt University Medical Center

#### Motivating Application: Error-Prone EHR data

- The Vanderbilt Comprehensive Care Clinic (VCCC) maintains an EHR database with  $\approx$  1300 people living with HIV receiving care from the VCCC
- Substantial error in key variables, including occurrences of AIDS-defining events (ADEs) and dates of antiretroviral therapy (ART) initiation
- Causal estimand: Average causal effect of beginning ART within 1 month of first visit on 3-year ADE risk
- Key problem: Both the outcome and treatment of interest are measured with error

#### **Two-Phase Sampling to Address Measurement Error**

Outcomes and treatments stored in EHR data are often measured with substantial error In many settings, it is possible to **validate** a random subset of error-prone observations





In practice, validated subjects are often selected according to a **sampling rule** which depends on all initially-observed data:  $X, Y^*$  and  $A^*$ 

**Goal**: Construct semiparametric efficient estimators of counterfactual means  $\mathbb{E}[Y(a)]$ 

**Challenge**: Validation datasets are typically <u>small</u> in practice – sources of finitesample instability can play a prominent role in estimation

#### Assumptions

- 1. Consistency: Y = AY(1) + (1 A)Y(0)
- 2. Treatment positivity:  $\mathbb{P}(A = 1 | \mathbf{X}) \in (0, 1)$
- 3. No unmeasured confounding:  $(Y(1), Y(0)) \perp A | \mathbf{X}|$
- 4. A and Y missing-at-random:  $(A, Y) \perp \mathbb{I} R | \mathbf{Z}$ , where  $\mathbf{Z} = (\mathbf{X}, A^*, Y^*)$
- 5. Validation positivity:  $\mathbb{P}(R = 1 | \mathbf{Z}) \in (0, 1)$

#### **Our Contributions**

- We present two asymptotically equivalent, semiparametric efficient one-step estimators
- We document unique sources of finite-sample instability faced by each estimator
- We present modifications to improve finite-sample behavior, and construct an ensemble estimator designed to prioritize finite-sample efficiency
- We developed the R package **drcmd**, which implements the Approach 2 estimator in general two-phase sampling and missing data settings

#### References

[1] James M Robins, Andrea Rotnitzky, and Lue Ping Zhao. Estimation of regression coefficients when some regressors are not always observed. Journal of the American statistical Association, 89(427):846–866, 1994.

# Efficient Estimation of Causal Effects Under Two-Phase Sampling with **Error-Prone Outcome and Treatment Measurements**

We connect two **asymptotically equivalent** approaches for constructing semiparametric efficient one-step estimators, and propose an ensemble estimator that optimizes finite sample efficiency.

# **Approach 1: Observed Data Distribution**

#### **High-level idea**: Follow the standard semiparametric statistics pipeline causal estimand $\mathbb{E}[Y(a)]$

Ma chawlunder Accumptione 1 5

$$\mathbb{E}[Y(a)] = \psi_{a,1} = \mathbb{E}\left[\frac{\mathbb{E}[Y_a(\mathbf{Z})]}{\mathbb{E}\{\lambda_a\}}\right]$$
  
where  $\lambda_a(\mathbf{Z}) = \mathbb{P}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}) = \mathbb{E}(A = a | \mathbf{Z}, R = 1)$  and  $\mu_a(\mathbf{Z}, R = 1)$  a

$$\hat{\psi}_{a,1}^{\mathrm{OS}} = \hat{\psi}_{a,1}^{\mathrm{PI}} + \frac{1}{n} \sum_{i=1}^{n} \widehat{\mathrm{EIF}}(\psi_{a,1}) + \frac{1}{n} \sum_{i=1}^{n} \widehat{\mathrm{EIF}}(\psi_{a$$

## **Approach 2: Complete Data Distribution**

**High-level idea**: [1] With complete data, could construct standard AIPW estimator

$$\hat{\psi}_{a,\mathsf{C}}^{\mathsf{Pl}} = \frac{1}{n} \sum_{i=1}^{n} \hat{m}_{a}(\boldsymbol{X}_{i}) \quad \text{and} \quad \hat{\psi}_{a,\mathsf{C}}^{\mathsf{OS}} = \hat{\psi}_{a,\mathsf{C}}^{\mathsf{Pl}} + \frac{1}{n} \sum_{i=1}^{n} \underbrace{\left(\hat{m}_{a}(\boldsymbol{X}_{i}) + \frac{I(A=a)}{\hat{g}_{a}(\boldsymbol{X}_{i})} \{Y - \hat{m}_{a}(\boldsymbol{X}_{i})\} - \hat{\psi}_{a,\mathsf{C}}^{\mathsf{Pl}}\right)}_{:=\boldsymbol{\chi}_{a}(\boldsymbol{O}; \hat{m}_{a}, \hat{g}_{a})}$$

where  $m_a(\mathbf{X}) = \mathbb{E}(Y|A = a, \mathbf{X})$  and  $g_a(\mathbf{X}) = \mathbb{P}(A = a|\mathbf{X})$  can be fit with weighted regressions that add weights  $R/\hat{\mathbb{P}}(R=1|\mathbf{Z})$  to the underlying loss functions.

Treating  $\chi_a(O; \hat{m}_a, \hat{g}_a)$  as pseudo-outcomes,

$$\hat{\psi}_{a,2}^{\text{OS}} = \hat{\psi}_{a,\text{C}}^{\text{PI}} + \frac{1}{n} \sum_{i=1}^{n} \left\{ \hat{\varphi}_{a}(\boldsymbol{Z}_{i}) + \frac{R_{i}}{\hat{\mathbb{P}}(R_{i}=1|\boldsymbol{Z}_{i})} \{ \boldsymbol{\chi}_{a}(\boldsymbol{O}_{i}; \hat{m}_{a}, \hat{g}_{a}) - \hat{\varphi}_{a}(\boldsymbol{Z}_{i}) \} \right\}$$

where  $\varphi_a(\mathbf{Z}) = \mathbb{E}[\chi_a(\mathbf{O}; m_a, g_a) | \mathbf{Z}, R = 1]$ , is an efficient one-step estimator.

#### **Properties**

Under Assumptions 1-5 and standard regularity conditions, we show that  $\hat{\psi}_{a,1}^{OS}$  and  $\hat{\psi}_{a,2}^{OS}$ are asymptotically equivalent and semiparametric efficient

Approach 2 can be viewed as a reparametrization of Approach 1

Both approaches have numerous unique sources of finite-sample instability

- Approach 1: debiasing term introduces numerous multiplicative, unstable weighting terms and requires estimation of 6 nuisance functions
- Approach 2: estimation of  $\varphi_a(\mathbf{Z})$  is an inherently difficult task in small samples

We propose estimating  $\varphi_a(Z)$  to minimize the **empirical variance** of  $\hat{\psi}_a^{OS,2}$ 

## **Ensemble Estimator**

**High-level idea**:  $\hat{\psi}_{a,1}^{OS}$  and  $\hat{\psi}_{a,2}^{OS}$  can differ substantially in finite samples  $\hat{\psi}_{a,\mathsf{E}}^{\mathsf{OS}} = \hat{w} \cdot \hat{\psi}_{a,1}^{\mathsf{OS}} + (1 - \hat{w}) \cdot \hat{\psi}_{a,2}^{\mathsf{OS}},$ 

where  $\hat{w}$  is chosen in a manner that (i) minimizes finite-sample variance, and (ii) remains well-defined asymptotically

debiased est.  $\hat{\psi}_a^{OS}$ 

,  $\lceil \mathbb{E} \{ \lambda_a(oldsymbol{Z}) \cdot \mu_a(oldsymbol{Z}) | oldsymbol{X} \} 
ceil$  $|(\boldsymbol{Z})|\boldsymbol{X}\}$  $(Y|A = a, \mathbf{Z}, R = 1)$  are imputation itor

 $\psi_{a,1}, O_i$ ).



## Real Data: Vanderbilt Comprehensive Care Clinic (VCCC)

Team of researchers validated **every** observation in VCCC database

 $\boldsymbol{Z}$ 



- large **bias** in estimate of ATE
- considered





# Rachel C. Nethery <sup>1</sup>

# Simulation

• Revealed increasingly larger shares of validated data, in manner which depends on

• High rates of measurement error in ADEs ( $\approx 12.5\%$ ) and early ART (4.1%) leads to

Ensemble estimator achieves lowest RMSE for all validation proportions