

# Shrinkage Estimation for Causal Inference and Experimental Design

**Evan Rosenman**<sup>†</sup>, Guillaume Basse, Mike Baiocchi,  
Art Owen, and Luke Miratrix

<sup>†</sup>Harvard University Data Science Initiative

January 18, 2023

# Outline

- 1 Problem Background
- 2 Assumptions and Set-Up
- 3 Inference
  - A Recipe for Estimators
  - Application to the WHI
- 4 Design
  - Problem Framework
  - Design Heuristics

# Motivating Setting

## Randomized Controlled Trials (RCT)

- Researcher controls assignment to treatment

## Observational Databases (ODB)

- Treatment assignments observed, but not controlled

# Motivating Setting

## Randomized Controlled Trials (RCT)

- Researcher controls assignment to treatment
  - Relatively few assumptions for unbiasedness
  - Often costly, small

## Observational Databases (ODB)

- Treatment assignments observed, but not controlled
  - Confounding  $\implies$  unverifiable assumptions for unbiasedness
  - Large, often inexpensive.

# Motivating Setting

## Randomized Controlled Trials (RCT)

- Researcher controls assignment to treatment
  - Relatively few assumptions for unbiasedness
  - Often costly, small
- “Unbiased but imprecise”

## Observational Databases (ODB)

- Treatment assignments observed, but not controlled
  - Confounding  $\implies$  unverifiable assumptions for unbiasedness
  - Large, often inexpensive.
- “Precise, but biased”

# Why should we care?

- Ubiquity of observational data in modern era
  - Electronic health records, disease surveillance
  - Fitness trackers, wearable devices, “internet of things”
  - E-commerce data, online behavior

# Why should we care?

- Ubiquity of observational data in modern era
  - Electronic health records, disease surveillance
  - Fitness trackers, wearable devices, “internet of things”
  - E-commerce data, online behavior
- Two major utilities to these data







# Our Approach

We consider two problems:

- **How to design shrinkage estimators to merge ODB and RCT data?**
- **How to improve experimental design using shrinkers?**

Work in a stratified setting, arising from:

- Subject matter knowledge
- Modern machine learning technique ([Wager and Athey, 2018](#); [Hill, 2011](#))

# Outline

- 1 Problem Background
- 2 Assumptions and Set-Up
- 3 Inference
  - A Recipe for Estimators
  - Application to the WHI
- 4 Design
  - Problem Framework
  - Design Heuristics

# Potential Outcomes Framework

- Have a sample of units  $i = 1, \dots, n$ . We are interested in some outcome measure  $Y$

# Potential Outcomes Framework

- Have a sample of units  $i = 1, \dots, n$ . We are interested in some outcome measure  $Y$
- For each unit,  $i$ , we suppose there are two associated values
  - $Y_i(1)$ : outcome if unit  $i$  receives the treatment
  - $Y_i(0)$ : outcome if unit  $i$  receives placebo

# Potential Outcomes Framework

- Have a sample of units  $i = 1, \dots, n$ . We are interested in some outcome measure  $Y$
- For each unit,  $i$ , we suppose there are two associated values
  - $Y_i(1)$ : outcome if unit  $i$  receives the treatment
  - $Y_i(0)$ : outcome if unit  $i$  receives placebo
- Causal quantity we are interested in is

$$\tau_i = Y_i(1) - Y_i(0)$$

# Causal Estimands

- **Fundamental Problem of Causal Inference**

- Each unit has a treatment status  $Z_i \in \{0, 1\}$ , and we observe

$$Y_i = Z_i Y_i(1) + (1 - Z_i) Y_i(0).$$

- Hence: cannot observe both  $Y_i(0)$  and  $Y_i(1)$  simultaneously!
- Typically settle for:
  - **Average treatment effect (ATE):**

$$\mathbb{E}(Y(1) - Y(0)), \quad \text{or}$$

- **Conditional average treatment effect (CATE):**

$$\mathbb{E}(Y(1) - Y(0) \mid X \in \mathcal{X}).$$

# Our Problem: Notation

- Observational data:  $n_o$  units sampled from

$$\left( \underbrace{Y_i(0), Y_i(1)}_{\text{potential outcomes}}, \underbrace{X_i}_{\text{covariates}}, \underbrace{Z_i}_{\text{treatment indicators}} \right) \stackrel{\text{iid}}{\sim} F_O.$$

- Experimental data: sample  $n_r$  units via

$$(Y_i(0), Y_i(1), X_i, Z_i) \stackrel{\text{iid}}{\sim} F_R.$$

- Assume strata  $k = 1, \dots, K$ . Stratum  $k$  defined by set of covariates values  $\mathcal{X}_k$ . Define indicators:

$$S_i = k \iff X_i \in \mathcal{X}_k.$$



# Assumptions and Non-Assumptions

- ① Under  $F_O$ ,

$$Y_i(1), Y_i(0) \mid X_i \underbrace{\quad}_{\text{"not independent of"}} Z_i$$

No unconfoundedness assumption for observational study.

# Assumptions and Non-Assumptions

- ① Under  $F_O$ ,

$$Y_i(1), Y_i(0) \mid X_i \not\perp Z_i$$

“not independent of”

No unconfoundedness assumption for observational study.

- ② Under  $F_R$ ,

$$Y_i(1), Y_i(0) \mid X_i \perp Z_i$$

“independent of”

# Assumptions and Non-Assumptions

- ① Under  $F_O$ ,

$$Y_i(1), Y_i(0) \mid X_i \not\perp Z_i$$

“not independent of”

No unconfoundedness assumption for observational study.

- ② Under  $F_R$ ,

$$Y_i(1), Y_i(0) \mid X_i \perp Z_i$$

“independent of”

- ③ For  $k = 1, \dots, K$ , have

$$\tau_k \equiv \mathbb{E}_R(Y_i(1) - Y_i(0) \mid S_i = k) = \mathbb{E}_O(Y_i(1) - Y_i(0) \mid S_i = k)$$

Assume **transportability** of CATEs across datasets.

Denote as  $\tau = (\tau_1, \dots, \tau_K)$  the vector of CATEs

# Setup

- Collect our estimators into vectors:

$$\hat{\boldsymbol{\tau}}_{\mathbf{r}} = (\hat{\tau}_{r1}, \dots, \hat{\tau}_{rK}), \quad \hat{\boldsymbol{\tau}}_{\mathbf{o}} = (\hat{\tau}_{o1}, \dots, \hat{\tau}_{oK}).$$

# Setup

- Collect our estimators into vectors:

$$\hat{\tau}_r = (\hat{\tau}_{r1}, \dots, \hat{\tau}_{rK}), \quad \hat{\tau}_o = (\hat{\tau}_{o1}, \dots, \hat{\tau}_{oK}).$$

- Under mild conditions, we have

$$\hat{\tau}_r \sim N(\tau, \Sigma_r), \quad \hat{\tau}_o \sim (\tau + \xi, \Sigma_o)$$

for bias  $\xi$  and diagonal covariance matrices  $\Sigma_r$  and  $\Sigma_o$

- $\Sigma_r = \text{diag}(\sigma_{r1}^2, \dots, \sigma_{rK}^2)$  is estimable from the data
- $\xi$  cannot be estimated

# Setup

- Collect our estimators into vectors:

$$\hat{\tau}_r = (\hat{\tau}_{r1}, \dots, \hat{\tau}_{rK}), \quad \hat{\tau}_o = (\hat{\tau}_{o1}, \dots, \hat{\tau}_{oK}).$$

- Under mild conditions, we have

$$\hat{\tau}_r \sim N(\tau, \Sigma_r), \quad \hat{\tau}_o \sim (\tau + \xi, \Sigma_o)$$

for bias  $\xi$  and diagonal covariance matrices  $\Sigma_r$  and  $\Sigma_o$

- $\Sigma_r = \text{diag}(\sigma_{r1}^2, \dots, \sigma_{rK}^2)$  is estimable from the data
- $\xi$  cannot be estimated
- Seek to design shrinkage estimator  $\hat{\tau} = f(\hat{\tau}_r, \hat{\tau}_o)$  to minimize expected squared error loss,

$$\mathcal{L}(\hat{\tau}, \tau) = \sum_k (\hat{\tau}_k - \tau_k)^2.$$

# Useful Prior Work

- **Shrinkage estimation**: “learn weights from the data”  $\implies$  a rich literature stretching back to multivariate normal mean estimation via the **James-Stein estimator** (Stein, 1956)

# Useful Prior Work

- **Shrinkage estimation**: “learn weights from the data”  $\implies$  a rich literature stretching back to multivariate normal mean estimation via the **James-Stein estimator** (Stein, 1956)
- Green and Strawderman (1991) and Green et al. (2005) propose estimators  $\delta_1, \delta_2$  for shrinkage between ...
  - a normal, unbiased estimator (like  $\hat{\tau}_r$ ), and
  - a biased estimator (like  $\hat{\tau}_o$ )



# Useful Prior Work

- **Shrinkage estimation:** “learn weights from the data”  $\implies$  a rich literature stretching back to multivariate normal mean estimation via the **James-Stein estimator** (Stein, 1956)
- Green and Strawderman (1991) and Green et al. (2005) propose estimators  $\delta_1, \delta_2$  for shrinkage between ...
  - a normal, unbiased estimator (like  $\hat{\tau}_r$ ), and
  - a biased estimator (like  $\hat{\tau}_o$ )
- **Key ideas**
  - Take weighted average of components of  $\hat{\tau}_r$  and  $\hat{\tau}_o$ .
  - Bias-variance tradeoff: estimators can stabilize high-variance  $\hat{\tau}_r$  by introducing some bias with shrinkage toward  $\hat{\tau}_o$

# Outline

- 1 Problem Background
- 2 Assumptions and Set-Up
- 3 Inference**
  - A Recipe for Estimators
  - Application to the WHI
- 4 Design
  - Problem Framework
  - Design Heuristics

# Outline

- 1 Problem Background
- 2 Assumptions and Set-Up
- 3 Inference
  - A Recipe for Estimators
  - Application to the WHI
- 4 Design
  - Problem Framework
  - Design Heuristics

# A Generalized Unbiased Risk Estimate (I)

## Theorem (Estimator Risk)

*Suppose we have  $\mathbf{U} \sim \mathcal{N}(\boldsymbol{\theta}, \boldsymbol{\Sigma})$ , random  $\mathbf{B}$ , and  $\mathcal{L}(\boldsymbol{\theta}, \mathbf{v}) = (\mathbf{v} - \boldsymbol{\theta})^\top \mathbf{W}(\mathbf{v} - \boldsymbol{\theta})$  where  $\boldsymbol{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_K^2)$  and  $\mathbf{W} = 1/K \cdot \text{diag}(w_1, \dots, w_K)$  is a diagonal weight matrix.*

# A Generalized Unbiased Risk Estimate (I)

## Theorem (Estimator Risk)

Suppose we have  $\mathbf{U} \sim \mathcal{N}(\boldsymbol{\theta}, \boldsymbol{\Sigma})$ , random  $\mathbf{B}$ , and  $\mathcal{L}(\boldsymbol{\theta}, \mathbf{v}) = (\mathbf{v} - \boldsymbol{\theta})^\top \mathbf{W}(\mathbf{v} - \boldsymbol{\theta})$  where  $\boldsymbol{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_K^2)$  and  $\mathbf{W} = 1/K \cdot \text{diag}(w_1, \dots, w_K)$  is a diagonal weight matrix. Then for

$$\kappa(\mathbf{U}, \mathbf{B}) = \mathbf{U} - \boldsymbol{\Sigma} \mathbf{g}(\mathbf{U}, \mathbf{B})$$

where  $\mathbf{g}(\mathbf{U}, \mathbf{B})$  is a function of  $\mathbf{U}$  and  $\mathbf{B}$  that is differentiable, satisfying  $E(\|\mathbf{g}\|^2) < \infty$ ,

# A Generalized Unbiased Risk Estimate (I)

## Theorem (Estimator Risk)

Suppose we have  $\mathbf{U} \sim \mathcal{N}(\boldsymbol{\theta}, \boldsymbol{\Sigma})$ , random  $\mathbf{B}$ , and  $\mathcal{L}(\boldsymbol{\theta}, \mathbf{v}) = (\mathbf{v} - \boldsymbol{\theta})^\top \mathbf{W}(\mathbf{v} - \boldsymbol{\theta})$  where  $\boldsymbol{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_K^2)$  and  $\mathbf{W} = 1/K \cdot \text{diag}(w_1, \dots, w_K)$  is a diagonal weight matrix. Then for

$$\kappa(\mathbf{U}, \mathbf{B}) = \mathbf{U} - \boldsymbol{\Sigma} \mathbf{g}(\mathbf{U}, \mathbf{B})$$

where  $\mathbf{g}(\mathbf{U}, \mathbf{B})$  is a function of  $\mathbf{U}$  and  $\mathbf{B}$  that is differentiable, satisfying  $E(\|\mathbf{g}\|^2) < \infty$ , we have

$$\begin{aligned} R(\boldsymbol{\theta}, \kappa(\mathbf{U}, \mathbf{B})) &= \mathbb{E}(\mathcal{L}(\boldsymbol{\theta}, \kappa(\mathbf{U}, \mathbf{B}))) \\ &= \frac{1}{K} \left( \text{Tr}(\boldsymbol{\Sigma} \mathbf{W}) + \mathbb{E} \left( \sum_{k=1}^K \sigma_k^4 w_k \left( g_k^2(\mathbf{U}, \mathbf{B}) - 2 \frac{\partial g_k(\mathbf{U}, \mathbf{B})}{\partial U_k} \right) \right) \right). \end{aligned}$$

# A Generalized Unbiased Risk Estimate (II)

From Theorem 1, obtain a generalization of Stein's Unbiased Risk Estimate ([Stein, 1981](#)),

$$\text{URE}(\boldsymbol{\theta}, \kappa(\mathbf{Z}, \mathbf{Y})) =$$

$$\frac{1}{K} \left( \text{Tr}(\boldsymbol{\Sigma} \mathbf{W}) + \sum_{k=1}^K \sigma_{rk}^4 w_k \left( g_k^2(\mathbf{U}, \mathbf{B}) - 2 \frac{\partial \mathbf{g}_k(\mathbf{U}, \mathbf{B})}{\partial U_k} \right) \right).$$

# A Generalized Unbiased Risk Estimate (II)

From Theorem 1, obtain a generalization of Stein's Unbiased Risk Estimate (Stein, 1981),

$$\text{URE}(\boldsymbol{\theta}, \kappa(\mathbf{Z}, \mathbf{Y})) = \frac{1}{K} \left( \text{Tr}(\boldsymbol{\Sigma} \mathbf{W}) + \sum_{k=1}^K \sigma_{rk}^4 w_k \left( g_k^2(\mathbf{U}, \mathbf{B}) - 2 \frac{\partial \mathbf{g}_k(\mathbf{U}, \mathbf{B})}{\partial U_k} \right) \right).$$

Common tactic: minimize URE over a hyperparameter (Li et al., 1985; Xie et al., 2012).



# A Generalized Unbiased Risk Estimate (II)

From Theorem 1, obtain a generalization of Stein's Unbiased Risk Estimate ([Stein, 1981](#)),

$$\text{URE}(\boldsymbol{\theta}, \boldsymbol{\kappa}(\mathbf{Z}, \mathbf{Y})) = \frac{1}{K} \left( \text{Tr}(\boldsymbol{\Sigma} \mathbf{W}) + \sum_{k=1}^K \sigma_{rk}^4 w_k \left( g_k^2(\mathbf{U}, \mathbf{B}) - 2 \frac{\partial \mathbf{g}_k(\mathbf{U}, \mathbf{B})}{\partial U_k} \right) \right).$$

Common tactic: minimize URE over a hyperparameter ([Li et al., 1985](#); [Xie et al., 2012](#)).

Points us toward a simple procedure:

- 1 Posit a structure for the shrinkage estimator
- 2 Derive a functional form by minimizing URE

# Case 1: Common Shrinkage Factor

We consider shrinkage estimators which share a common shrinkage  $\lambda$  factor across components. Denote a generic estimator as

$$\kappa(\lambda, \hat{\tau}_r, \hat{\tau}_o) = \hat{\tau}_r - \lambda(\hat{\tau}_r - \hat{\tau}_o).$$

# Case 1: Common Shrinkage Factor

We consider shrinkage estimators which share a common shrinkage  $\lambda$  factor across components. Denote a generic estimator as

$$\kappa(\lambda, \hat{\tau}_r, \hat{\tau}_o) = \hat{\tau}_r - \lambda(\hat{\tau}_r - \hat{\tau}_o).$$

Then the URE evaluates to

$$\text{URE}(\lambda) = \text{Tr}(\Sigma_r \mathbf{W}) + \lambda^2 (\hat{\tau}_o - \hat{\tau}_r)^T \mathbf{W} (\hat{\tau}_o - \hat{\tau}_r) - 2\lambda \text{Tr}(\Sigma_r \mathbf{W})$$

# Case 1: Common Shrinkage Factor

We consider shrinkage estimators which share a common shrinkage  $\lambda$  factor across components. Denote a generic estimator as

$$\kappa(\lambda, \hat{\tau}_r, \hat{\tau}_o) = \hat{\tau}_r - \lambda(\hat{\tau}_r - \hat{\tau}_o).$$

Then the URE evaluates to

$$\text{URE}(\lambda) = \text{Tr}(\Sigma_r \mathbf{W}) + \lambda^2 (\hat{\tau}_o - \hat{\tau}_r)^\top \mathbf{W} (\hat{\tau}_o - \hat{\tau}_r) - 2\lambda \text{Tr}(\Sigma_r \mathbf{W})$$

which has minimizer in  $\lambda$ ,

$$\lambda_1^{\text{URE}} = \frac{\text{Tr}(\Sigma_r \mathbf{W})}{(\hat{\tau}_o - \hat{\tau}_r)^\top \mathbf{W} (\hat{\tau}_o - \hat{\tau}_r)}.$$

# A Note on $\lambda_1^{\text{URE}}$

The true risk-minimizing shrinkage weight is given by

$$\lambda_{\text{opt}} = \frac{\text{Tr}(\Sigma_r \mathbf{W})}{\text{Tr}(\Sigma_r \mathbf{W}) + \text{Tr}(\Sigma_o \mathbf{W}) + \underbrace{\xi^\top \mathbf{W} \xi}_{\text{Not estimable from data}}},$$

but observe that

$$E \left( (\hat{\tau}_o - \hat{\tau}_r)^\top \mathbf{W} (\hat{\tau}_o - \hat{\tau}_r) \right) = \text{Tr}(\Sigma_r \mathbf{W}) + \text{Tr}(\Sigma_o \mathbf{W}) + \xi^\top \mathbf{W} \xi.$$

$\lambda_1^{\text{URE}}$  substitutes the quadratic form for its expectation,

$$\lambda_1^{\text{URE}} = \frac{\text{Tr}(\Sigma_r \mathbf{W})}{(\hat{\tau}_o - \hat{\tau}_r)^\top \mathbf{W} (\hat{\tau}_o - \hat{\tau}_r)}.$$

# Useful Properties of $\lambda_1^{\text{URE}}$ (I)

## 1 Define

$$\kappa_1 = \hat{\tau}_r - \lambda_1^{\text{URE}} (\hat{\tau}_r - \hat{\tau}_o)$$

$\kappa_1$  admits a testable condition under which it is guaranteed to reduce risk relative to  $\hat{\tau}_r$ .

### Lemma ( $\kappa_1$ Risk Guarantee)

*Suppose  $4 \max_k w_k \sigma_{rk}^2 < \sum_k w_k \sigma_{rk}^2$ . Then  $\kappa_1$  has risk strictly less than that of  $\hat{\tau}_r$ .*

- Requires a dimension of at least  $K = 5$ .
- May require substantially larger  $K$  if high heteroscedasticity or non-uniform weights.

# Useful Properties of $\lambda_1^{\text{URE}}$ (II)

- ② Its positive part analogue,

$$\kappa_{1+} = \hat{\tau}_r - \left\{ \lambda_1^{\text{URE}} \right\}_{[0,1]} (\hat{\tau}_r - \hat{\tau}_o) ,$$

where

$$\{u\}_{[0,1]} = \min(\max(u, 0), 1) ,$$

satisfies the following notion of optimality:

# Useful Properties of $\lambda_1^{\text{URE}}$ (III)

## Theorem ( $\kappa_{1+}$ Asymptotic Risk)

*Suppose*

$$\limsup_{K \rightarrow \infty} \frac{1}{K} \sum_k d_k^2 \sigma_{rk}^2 \xi_k^2 < \infty, \quad \limsup_{K \rightarrow \infty} \frac{1}{K} \sum_k d_k^2 \sigma_{rk}^2 \sigma_{ok}^2 < \infty,$$

$$\text{and} \quad \limsup_{K \rightarrow \infty} \frac{1}{K} \sum_k d_k^2 \sigma_{rk}^4 < \infty.$$

*Then, in the limit  $K \rightarrow \infty$ ,  $\kappa_{1+}$  has the lowest risk among all estimators with a shared shrinkage factor across components.*



## Case 2: Variance-Weighted Shrinkage Factor

This procedure is general purpose. For example, may instead want an estimator that shrinks each component proportionally to  $\sigma_{rk}^2$ .

Easy to solve for

$$\kappa_2 = \kappa(\lambda_2^{\text{URE}}, \hat{\tau}_r, \hat{\tau}_o) = \hat{\tau}_r - \frac{\text{Tr}(\Sigma_r^2 \mathbf{W}) \Sigma_r}{(\hat{\tau}_o - \hat{\tau}_r)^\top \Sigma_r^2 \mathbf{W} (\hat{\tau}_o - \hat{\tau}_r)} (\hat{\tau}_r - \hat{\tau}_o)$$

and its positive-part improvement,

$$\kappa_{2+} = \hat{\tau}_r - \left\{ \frac{\text{Tr}(\Sigma_r^2 \mathbf{W}) \Sigma_r}{(\hat{\tau}_o - \hat{\tau}_r)^\top \Sigma_r^2 \mathbf{W} (\hat{\tau}_o - \hat{\tau}_r)} \right\}_{[0,1]} (\hat{\tau}_r - \hat{\tau}_o) .$$



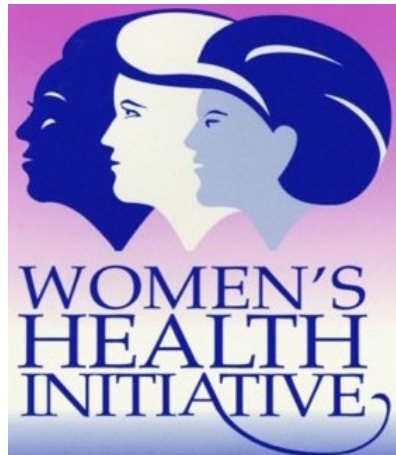
# Outline

- 1 Problem Background
- 2 Assumptions and Set-Up
- 3 Inference
  - A Recipe for Estimators
  - Application to the WHI
- 4 Design
  - Problem Framework
  - Design Heuristics

# WHI Overview

## Dataset Overview

- Study of postmenopausal women initiated in 1991
- RCT of hormone therapy (estrogen and progestin) w/ 16k enrollees
- ODB w/ 50k comparable enrollees



# Application to the WHI

- Compute “true” causal effect of **hormone therapy** on **coronary heart disease** using entire RCT (16k units)

# Application to the WHI

- Compute “true” causal effect of **hormone therapy** on **coronary heart disease** using entire RCT (16k units)
- Repeat 500 times:
  - Draw bootstrap samples:
    - 1,000 RCT units
    - Observational sample (50k units)
  - Compute squared error loss for  $\hat{\tau}_r, \kappa_{1+}, \kappa_{2+}, \delta_1, \delta_2$ .

# Application to the WHI

- Compute “true” causal effect of **hormone therapy** on **coronary heart disease** using entire RCT (16k units)
- Repeat 500 times:
  - Draw bootstrap samples:
    - 1,000 RCT units
    - Observational sample (50k units)
  - Compute squared error loss for  $\hat{\tau}_r, \kappa_{1+}, \kappa_{2+}, \delta_1, \delta_2$ .
- Average loss over draws

# Choice of Stratification Variables

Stratify on:

- two variables from WHI protocol:  
age + history of cardiovascular disease (Roehm, 2015).
- a variable unassociated with treatment effect:  
solar irradiance (“Langley”)  $\implies$  uncorrelated w/ outcome



## Results

Subgroup Variable(s)	# of Strata	Loss as % of $\hat{\tau}_r$ Loss			
		$\kappa_{1+}$	$\kappa_{2+}$	$\delta_1$	$\delta_2$
CVD	2	37.6%	<b>36.9%</b>	100.0%	100.0%
Age	3	37.3%	<b>30.1%</b>	61.5%	72.8%
Langley	5	29.4%	<b>23.5%</b>	40.0%	52.2%
CVD, Age	6	<b>38.0%</b>	38.2%	38.3%	82.4%
CVD, Langley	10	30.6%	32.5%	<b>30.0%</b>	87.2%
Age, Langley	15	<b>22.4%</b>	23.0%	22.5%	43.1%
Age, CVD, Langley	30	50.3%	<b>50.3%</b>	50.3%	78.4%

# Outline

- 1 Problem Background
- 2 Assumptions and Set-Up
- 3 Inference
  - A Recipe for Estimators
  - Application to the WHI
- 4 Design
  - Problem Framework
  - Design Heuristics



# A New Setting: Design

Can these insights inform the design of a **prospective** RCT?

- Observational study already completed,  $\hat{\tau}_o$  obtained.
- Designing a prospective RCT of  $n_r$  units
- Want to use a shrinker to combine  $\hat{\tau}_r$  with  $\hat{\tau}_o$ . Design experiment to better complement ODB

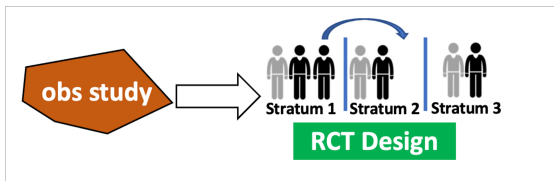
# A New Setting: Design

Can these insights inform the design of a **prospective** RCT?

- Observational study already completed,  $\hat{\tau}_o$  obtained.
- Designing a prospective RCT of  $n_r$  units
- Want to use a shrinker to combine  $\hat{\tau}_r$  with  $\hat{\tau}_o$ . Design experiment to better complement ODB

**Goal:** choose an RCT allocation of treated and control counts per stratum,  $\mathbf{d} = \{(n_{rkt}, n_{rkc})\}_{k=1}^K$ , s.t.  $\sum_k n_{rkt} + n_{rkc} = n_r$ :

- implies how to *recruit* ...
- and *assign* treatment



# Estimator and Risk

We proceed with our estimator  $\kappa_{2+}$  from the prior section:

$$\kappa_{2+} = \hat{\tau}_r - \left\{ \frac{\text{Tr}(\Sigma_r^2 \mathbf{W}) \Sigma_r}{(\hat{\tau}_o - \hat{\tau}_r)^\top \Sigma_r^2 \mathbf{W} (\hat{\tau}_o - \hat{\tau}_r)} \right\}_{[0,1]} (\hat{\tau}_r - \hat{\tau}_o)$$

# Estimator and Risk

We proceed with our estimator  $\kappa_{2+}$  from the prior section:

$$\kappa_{2+} = \hat{\tau}_r - \left\{ \frac{\text{Tr}(\Sigma_r^2 \mathbf{W}) \Sigma_r}{(\hat{\tau}_o - \hat{\tau}_r)^\top \Sigma_r^2 \mathbf{W} (\hat{\tau}_o - \hat{\tau}_r)} \right\}_{[0,1]} (\hat{\tau}_r - \hat{\tau}_o)$$

Optimize experimental design over  $\mathcal{R}_2(\mathbf{d}, \mathbf{V}, \xi)$ , the risk of  $\kappa_{2+}$  under fixed  $\hat{\tau}_o$ , with

# Estimator and Risk

We proceed with our estimator  $\kappa_{2+}$  from the prior section:

$$\kappa_{2+} = \hat{\tau}_r - \left\{ \frac{\text{Tr}(\Sigma_r^2 \mathbf{W}) \Sigma_r}{(\hat{\tau}_o - \hat{\tau}_r)^\top \Sigma_r^2 \mathbf{W} (\hat{\tau}_o - \hat{\tau}_r)} \right\}_{[0,1]} (\hat{\tau}_r - \hat{\tau}_o)$$

Optimize experimental design over  $\mathcal{R}_2(\mathbf{d}, \mathbf{V}, \xi)$ , the risk of  $\kappa_{2+}$  under fixed  $\hat{\tau}_o$ , with

- design  $\mathbf{d}$



# Estimator and Risk

We proceed with our estimator  $\kappa_{2+}$  from the prior section:

$$\kappa_{2+} = \hat{\tau}_r - \left\{ \frac{\text{Tr}(\Sigma_r^2 \mathbf{W}) \Sigma_r}{(\hat{\tau}_o - \hat{\tau}_r)^\top \Sigma_r^2 \mathbf{W} (\hat{\tau}_o - \hat{\tau}_r)} \right\}_{[0,1]} (\hat{\tau}_r - \hat{\tau}_o)$$

Optimize experimental design over  $\mathcal{R}_2(\mathbf{d}, \mathbf{V}, \xi)$ , the risk of  $\kappa_{2+}$  under fixed  $\hat{\tau}_o$ , with

- design  $\mathbf{d}$
- stratum potential outcome variances  $\mathbf{V} = \{(\hat{\sigma}_{kt}^2, \hat{\sigma}_{kc}^2)\}_{k=1}^K$

# Estimator and Risk

We proceed with our estimator  $\kappa_{2+}$  from the prior section:

$$\kappa_{2+} = \hat{\tau}_r - \left\{ \frac{\text{Tr}(\Sigma_r^2 \mathbf{W}) \Sigma_r}{(\hat{\tau}_o - \hat{\tau}_r)^\top \Sigma_r^2 \mathbf{W} (\hat{\tau}_o - \hat{\tau}_r)} \right\}_{[0,1]} (\hat{\tau}_r - \hat{\tau}_o)$$

Optimize experimental design over  $\mathcal{R}_2(\mathbf{d}, \mathbf{V}, \xi)$ , the risk of  $\kappa_{2+}$  under fixed  $\hat{\tau}_o$ , with

- design  $\mathbf{d}$
- stratum potential outcome variances  $\mathbf{V} = \{(\hat{\sigma}_{kt}^2, \hat{\sigma}_{kc}^2)\}_{k=1}^K$
- bias vector  $\xi$ .

# Estimator and Risk

We proceed with our estimator  $\kappa_{2+}$  from the prior section:

$$\kappa_{2+} = \hat{\tau}_r - \left\{ \frac{\text{Tr}(\Sigma_r^2 \mathbf{W}) \Sigma_r}{(\hat{\tau}_o - \hat{\tau}_r)^\top \Sigma_r^2 \mathbf{W} (\hat{\tau}_o - \hat{\tau}_r)} \right\}_{[0,1]} (\hat{\tau}_r - \hat{\tau}_o)$$

Optimize experimental design over  $\mathcal{R}_2(\mathbf{d}, \mathbf{V}, \xi)$ , the risk of  $\kappa_{2+}$  under fixed  $\hat{\tau}_o$ , with

- design  $\mathbf{d}$
- stratum potential outcome variances  $\mathbf{V} = \{(\hat{\sigma}_{kt}^2, \hat{\sigma}_{kc}^2)\}_{k=1}^K$
- bias vector  $\xi$ .

Can compute this efficiently via numerical integration ([Bao and Kan, 2013](#)), as long as  $\mathbf{V}$  and  $\xi$  are known.

# Outline

- 1 Problem Background
- 2 Assumptions and Set-Up
- 3 Inference
  - A Recipe for Estimators
  - Application to the WHI
- 4 Design
  - Problem Framework
  - Design Heuristics

# 1. Neyman Allocation

Can estimate  $\mathbf{V}$  using pilot estimates obtained from ODB:

$$\hat{\sigma}_{kt}^2 = \widehat{\text{var}}(Y(1) \mid S = k) \quad \text{and} \quad \hat{\sigma}_{kc}^2 = \widehat{\text{var}}(Y(0) \mid S = k) .$$

Simplest design heuristic: use a Neyman allocation, e.g.

$$n_{rkt} = \frac{n_r \cdot \hat{\sigma}_{kt}^2}{\sum_k \hat{\sigma}_{kt}^2 + \hat{\sigma}_{kc}^2} \quad \text{and} \quad n_{rkc} = \frac{n_r \cdot \hat{\sigma}_{kc}^2}{\sum_k \hat{\sigma}_{kt}^2 + \hat{\sigma}_{kc}^2} .$$

Optimizes over only the non-shrinkage portion of the risk, but reasonable in many practical settings.

## 2. Naïve Optimization Assuming $\xi = 0$ (I)

Use. a simple heuristic: assume  $\xi = 0$ . Then solve:

$$\begin{aligned} & \text{minimize} && \mathcal{R}_2(\mathbf{d}, \mathbf{V}, \xi) \\ & \text{subject to} && \xi = 0, \mathbf{V} = \{(\hat{\sigma}_{kt}^2, \hat{\sigma}_{kc}^2)\}_{k=1}^K, \\ & && 0 < n_{rkt}, n_{rkc},, \quad k = 1, \dots, K, \\ & && n_r = \sum_k n_{rkt} + n_{rkc}. \end{aligned} \tag{1}$$

But  $\mathcal{R}_2(\mathbf{d}, \mathbf{V}, \xi)$  is not convex in the design  $\mathbf{d}$ ...

## 2. Naïve Optimization Assuming $\xi = 0$ (II)

A practical approach: **greedy algorithm**. Define  $\mathbf{d}_j$  as design on  $j^{th}$  iteration, and define

$$\mathcal{D}_j = \{\mathbf{d}' \mid \mathbf{d}' \text{ changes one unit across strata/treatment level from } \mathbf{d}_j\}.$$

Run Algorithm 2 from several values of  $\mathbf{d}_0$  and take minimum:

Start with design  $\mathbf{d}_0 = \{(n_{rkt}^{(0)}, n_{rkC}^{(0)})\}_k$ .

For iteration  $j = 1, 2, \dots$ :

For each design  $\mathbf{d}'$  in  $\mathcal{D}_{j-1}$ :

Compute  $\mathcal{R}_2(\mathbf{d}', \mathbf{V}, 0)$ .

(2)

Set  $\mathbf{d}_j = \underset{\mathbf{d}' \in \mathcal{D}_{j-1}}{\operatorname{argmin}} \mathcal{R}_2(\mathbf{d}', \mathbf{V}, 0)$

If  $\mathcal{R}_2(\mathbf{d}_j, \mathbf{V}, 0) \geq \mathcal{R}_2(\mathbf{d}_{j-1}, \mathbf{V}, 0)$

Return  $\mathbf{d}_{j-1}$ .

### 3. Heuristic Optimization Assuming Worst-Case Error Under $\Gamma$ -Level Unmeasured Confounding

- Can take a more pessimistic approach again using marginal sensitivity model of [Tan \(2006\)](#)
- For a user-chosen value of  $\Gamma \geq 1$ :
  - can obtain worst-case  $\xi_k(\Gamma)$  using [Zhao et al. \(2019\)](#), and...
  - if outcome  $Y_i \in \{0, 1\}$ , can obtain associated  $\hat{\sigma}_{kt}^2$  and  $\hat{\sigma}_{kc}^2$ .



### 3. Heuristic Optimization Assuming Worst-Case Error Under $\Gamma$ -Level Unmeasured Confounding

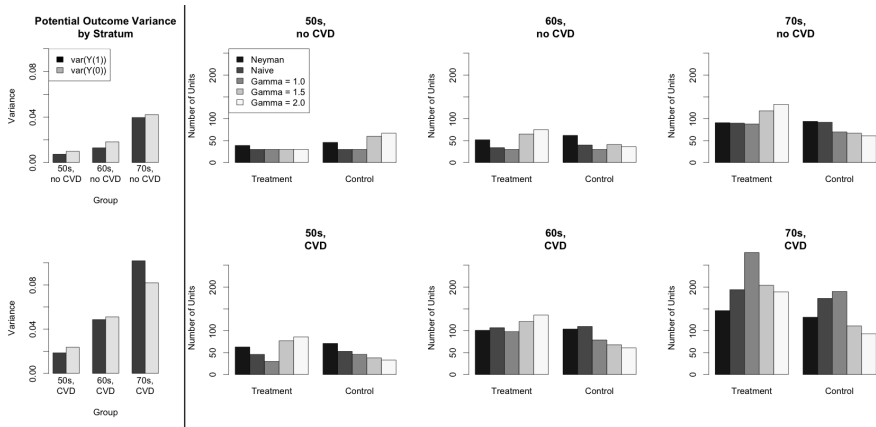
- Can take a more pessimistic approach again using marginal sensitivity model of [Tan \(2006\)](#)
- For a user-chosen value of  $\Gamma \geq 1$ :
  - can obtain worst-case  $\xi_k(\Gamma)$  using [Zhao et al. \(2019\)](#), and...
  - if outcome  $Y_i \in \{0, 1\}$ , can obtain associated  $\hat{\sigma}_{kt}^2$  and  $\hat{\sigma}_{kc}^2$ .

posit a value of  $\Gamma \implies$

collect results into  $\mathbf{V}(\Gamma)$  and  $\boldsymbol{\xi}(\Gamma) \implies$

run Algorithm 2 using  $\mathcal{R}_2(\mathbf{d}, \mathbf{V}(\Gamma), \boldsymbol{\xi}(\Gamma))$  instead

# Stratified WHI Study Design of $n_r = 1,000$ units



**Figure 2:** Allocations in WHI with strata defined by history of CVD and age, under different design heuristics.

## Future Work

Some areas I'm excited about pursuing:

- **Applied project:** air pollution and mortality (with Francesca Dominici & Luke Miratrix)
  - Combining Medicare (“observational database”) database with Medicare Current Beneficiary Survey (“close to” RCT)
  - Approach via *double shrinkage*:

$$\psi_k = a_k (\lambda_k \hat{\tau}_{rk} + (1 - \lambda_k) \hat{\tau}_{ok})$$

where  $a_k, \lambda_k$  are data-driven EB shrinkage parameters

- **ML approaches**
  - Move beyond stratification
  - Flexible shrinkage between CATE functions  $\hat{\tau}_r(x)$  and  $\hat{\tau}_o(x)$

# Acknowledgments

Thank you to my collaborators on this work:

- Guillaume Basse
- Art Owen
- Mike Baiocchi
- Luke Miratrix

Inference paper out in *Biometrics*.

Design paper available at [arXiv:2204.06687](https://arxiv.org/abs/2204.06687)











# Practical Considerations

- **Variance estimation:** In practice,  $\Sigma_r$  not known. Must be estimated from data.

# Practical Considerations

- **Variance estimation:** In practice,  $\Sigma_r$  not known. Must be estimated from data.
- **Propensity score adjustment**
  - No unconfoundedness  $\implies$   
propensity score adjustment can't remove all bias
  - If ODB is large, adjusting will typically be good practice. We suggest stabilized IPTW adjustments.



# Computing Shrinker Risk

Goal is to optimize experimental design over  $\mathcal{R}(\kappa_2)$ .

Define  $\mathcal{R}_2(\mathbf{d}, \mathbf{V}, \boldsymbol{\xi})$  as risk of  $\kappa_2$  under fixed  $\hat{\tau}_o$ , with

- design  $\mathbf{d}$
- stratum potential outcome variances  $\mathbf{V} = \{(\hat{\sigma}_{kt}^2, \hat{\sigma}_{kc}^2)\}_{k=1}^K$
- bias vector  $\boldsymbol{\xi}$ .

# Computing Shrinker Risk

Goal is to optimize experimental design over  $\mathcal{R}(\kappa_2)$ .

Define  $\mathcal{R}_2(\mathbf{d}, \mathbf{V}, \xi)$  as risk of  $\kappa_2$  under fixed  $\hat{\tau}_o$ , with

- design  $\mathbf{d}$
- stratum potential outcome variances  $\mathbf{V} = \{(\hat{\sigma}_{kt}^2, \hat{\sigma}_{kc}^2)\}_{k=1}^K$
- bias vector  $\xi$ .

Reduces to a ratio of Gaussian quadratic forms!  $\implies$   
solvable via numerical integral of [Bao and Kan \(2013\)](#)

**Upshot:** can efficiently compute the risk of any design if we have values for  $\mathbf{V}$  and  $\xi$ .

# Estimating $V$ : Updated Assumptions

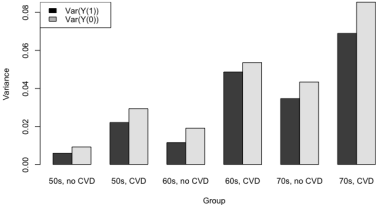
Same assumptions, but a stronger form of **transportability**:

- ③ For  $k = 1, \dots, K$  and  $w \in \{0, 1\}$ :

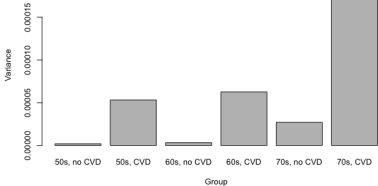
$$\mathbb{E}_O(Y(w) \mid S = k) = \mathbb{E}_R(Y(w) \mid S = k) \text{ and} \\ \text{var}_O(Y(w) \mid S = k) = \text{var}_R(Y(w) \mid S = k) .$$

# Sample Designs

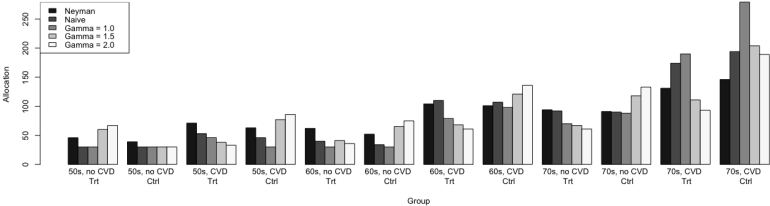
Potential Outcome Variance by Stratum



Observational Estimator Variance by Stratum



Allocations by Stratum Under Different Schemes



# Guardrails

Simplicity of Algorithm 2 makes it easy to impose guardrails  $\implies$   
for any invalid design, just set objective value to  $\infty$ .

Recommend simple guardrails for designs:

- 1 **Sample size:** to retain CLT, enforce

$$\min_k n_{rkt} \geq SS_{\min}, \quad \min_k n_{rk c} \geq SS_{\min}$$

- 2 **Detachability:** for default design  $\tilde{\mathbf{d}} = \{\tilde{n}_{rkt}, \tilde{n}_{rk c}\}_k$  and tolerance parameter  $\delta_d \geq 1$ , enforce

$$\sum_k \frac{\hat{\sigma}_{kt}^2}{n'_{rkt}} + \frac{\hat{\sigma}_{kc}^2}{n'_{rk c}} \geq \delta_d \sum_k \frac{\hat{\sigma}_{kt}^2}{\tilde{n}_{rkt}} + \frac{\hat{\sigma}_{kc}^2}{\tilde{n}_{rk c}},$$

for any proposed design  $\mathbf{d}' = \{n'_{rkt}, n'_{rk c}\}_k$ .



- Rich data set. Consider 684 covariates: demographics, medical history, diet, etc.
- Fit  $\hat{e}(\mathbf{x}) = \hat{\mathbb{E}}(W \mid \mathbf{x})$  by stepwise logistic regression w/ cross-validation. 53 variables chosen.



# Covariate Balance (I)

**Table 1:** Standardized differences (SD) between treated and control populations in the observational dataset, before and after stratification on the propensity score, for clinical risk factors for coronary heart disease.

	Unweighted			Stratified		
	Test	Ctrl	SD	Test	Ctrl	SD
<b>Age</b>	60.78	64.72	-0.56	63.06	63.33	-0.04
<b>BMI</b>	25.55	27.11	-0.25	26.71	26.62	0.00
<b>Physical functioning</b>	85.23	79.58	0.26	81.15	81.23	0.03
<b>Age at menopause</b>	50.49	50.19	0.06	50.35	50.33	0.02

# Covariate Balance (II)

**Table 2:** Standardized differences (SD) between treated and control populations in the observational database, before and after stratification on the propensity score, for ethnicity category.

		White	Black	Latino	AAPI	Native American	Missing/ Other	SD
<b>Before Strat.</b>	Treated	89.0%	2.7%	2.9%	4.0%	0.2%	1.1%	0.26
	Control	83.1%	8.1%	3.9%	2.8%	0.4%	1.5%	
<b>After Strat.</b>	Treated	83.4%	6.9%	4.3%	3.6%	0.5%	1.4%	0.05
	Control	84.8%	6.4%	3.6%	3.4%	0.4%	1.4%	

# Covariate Balance (III)

**Table 3:** Standardized differences (SD) between treated and control populations in the observational database, before and after stratification on the propensity score, for smoking category.

		<b>Never Smoked</b>	<b>Past Smoker</b>	<b>Current Smoker</b>	<b>SD</b>
<b>Before Stratifying</b>	Treated	48.7%	46.2%	5.1%	0.11
	Control	52.3%	41.1%	6.6%	
<b>After Stratifying</b>	Treated	50.9%	42.5%	6.6%	0.01
	Control	51.0%	42.7%	6.3%	

## Useful Prior Results (I)

- **Green and Strawderman (1991)** consider the  $\Sigma_o = \gamma^2 I_K, \Sigma_r = \sigma_r^2 I_K$  case. Show that estimator

$$\hat{\tau}_{\mathbf{o}} + \left(1 - \frac{(K-2)\sigma_r^2}{\|\hat{\tau}_{\mathbf{r}} - \hat{\tau}_{\mathbf{o}}\|^2}\right)_+ (\hat{\tau}_{\mathbf{r}} - \hat{\tau}_{\mathbf{o}})$$

dominates  $\hat{\tau}_r$  under squared error loss, and has bounded risk as  $\xi$  grows large

## Useful Prior Results (II)

- **Green et al. (2005)**: Generalize results to heteroscedastic case and propose modified estimators

$$\delta_1 = \hat{\tau}_o + \left( 1 - \frac{(K-2)}{(\hat{\tau}_r - \hat{\tau}_o)^\top \Sigma_r^{-1} (\hat{\tau}_r - \hat{\tau}_o)} \right)_+ (\hat{\tau}_r - \hat{\tau}_o)$$

$$\delta_2 = \hat{\tau}_o + \left( 1 - \frac{(K-2)\Sigma_r^{-1}}{(\hat{\tau}_r - \hat{\tau}_o)^\top \Sigma_r^{-2} (\hat{\tau}_r - \hat{\tau}_o)} \right)_+ (\hat{\tau}_r - \hat{\tau}_o)$$

Fewer theoretical guarantees.

$\delta_1$  is designed for precision-weighted loss, but outperforms  $\delta_2$  under regular  $L_2$  loss in simulation.

# Integral Expressions

Bao and Kan (2013) give a method for computing these ratios exactly via numerical integrals:

$$\mathbb{E}_r \left( \frac{\boldsymbol{\nu}^\top \boldsymbol{\Sigma}_r^5 \boldsymbol{\nu}}{(\boldsymbol{\nu}^\top \boldsymbol{\Sigma}_r^3 \boldsymbol{\nu})^2} \right) = \int_0^\infty \det(\mathbf{I} + 2t\boldsymbol{\Sigma}_r^3)^{-1/2} \cdot \exp \left( \frac{1}{2} (\boldsymbol{\xi}^\top (\mathbf{I} + 2t\boldsymbol{\Sigma}_r^3)^{-1} \boldsymbol{\xi} - \boldsymbol{\xi}^\top \boldsymbol{\xi}) \right) \\ \left( \text{Tr}(\mathbf{R}) + (\mathbf{L}\boldsymbol{\Sigma}_r^{-1/2}\boldsymbol{\xi})^\top \mathbf{R} (\mathbf{L}\boldsymbol{\Sigma}_r^{-1/2}\boldsymbol{\xi}) \right) t dt$$

$$\mathbb{E}_r \left( \frac{1}{(\boldsymbol{\nu}^\top \boldsymbol{\Sigma}_r^3 \boldsymbol{\nu})} \right) = \int_0^\infty \det(\mathbf{I} + 2t\boldsymbol{\Sigma}_r^3)^{-1/2} \cdot \exp \left( \frac{1}{2} (\boldsymbol{\xi}^\top (\mathbf{I} + 2t\boldsymbol{\Sigma}_r^3)^{-1} \boldsymbol{\xi} - \boldsymbol{\xi}^\top \boldsymbol{\xi}) \right) t dt$$

$$\text{where } \mathbf{L} = (\mathbf{I} + 2t\boldsymbol{\Sigma}_r^3)^{-1/2} \quad \text{and} \quad \mathbf{R} = \mathbf{L}^\top \boldsymbol{\Sigma}_r^5 \mathbf{L}.$$

This gives us a way to efficiently compute the risk of any design, under a set of assumptions about the values of  $\boldsymbol{\Sigma}_r$  and  $\boldsymbol{\xi}$ .

# Improving Interpretability of $\kappa_{1+}$

- Recall:  $\lambda_1^{\text{URE}}$  can be interpreted as an estimate of

$$\lambda_{\text{opt}} = \frac{\text{Tr}(\Sigma_r \mathbf{W})}{\text{Tr}(\Sigma_r \mathbf{W}) + \text{Tr}(\Sigma_o \mathbf{W}) + \xi^T \mathbf{W}^2 \xi},$$

true MSE-minimizing weight on  $\hat{\tau}_o$  in a convex combination

- We can use this idea to improve interpretability of  $\kappa_{1+}$ !
- Key idea:** frame in context of sensitivity model of [Tan \(2006\)](#)



# Prior Work

- Marginal sensitivity model of [Tan \(2006\)](#) summarizes degree of unmeasured confounding by a single value,  $\Gamma \geq 1$ 
  - $\Gamma$  bounds odds ratio of treatment prob. conditional on potential outcomes + covariates vs. covariates only
  - Related to the famous model of [Rosenbaum \(1987\)](#), but extends to the setting of inverse probability weighting
- [Zhao et al. \(2019\)](#) derive valid confidence intervals for causal estimates under the set of models indexed by any choice of  $\Gamma$ 
  - Implicitly maps  $\Gamma$  to a worst-case bias  $\xi(\Gamma)$  and variance  $\Sigma_O(\Gamma)$
  - Under some assumptions, allows us to obtain worst-case estimate of  $\lambda_{\text{opt}}$  as a function of  $\Gamma$ , which we call  $\lambda(\Gamma)$

# Relating the Models

- **Intuition:** larger  $\Gamma$  (confounding parameter)  $\implies$  optimal weight  $\lambda_{\text{opt}}$  is smaller
- Let  $\Gamma_{\text{imp}} = \sup\{\Gamma : \lambda(\Gamma) > \lambda_1^{\text{URE}}\}$ 
  - Largest value  $\Gamma$  for which the optimal shrinkage factor  $\lambda(\Gamma)$  is greater than our shrinkage parameter  $\lambda_1^{\text{URE}}$ .
- $\Gamma_{\text{imp}}$  can be used to evaluate level of shrinkage
  - If we believe true confounding level  $\Gamma < \Gamma_{\text{imp}}$ , then

$$\lambda_1^{\text{URE}} \approx \lambda(\Gamma_{\text{imp}}) \leq \lambda_{\text{opt}} = \lambda(\Gamma)$$

Hence the shrinkage level is conservative. ✓

- If we believe  $\Gamma > \Gamma_{\text{imp}}$ , then estimator is overshrinking, relies too much on the observational estimate. ✗

# Simulations Set-Up (I)

- ODB has 20K units ( $j \in \mathcal{O}$ ). RCT has 1,000 ( $i \in \mathcal{E}$ )
- Untreated potential outcomes  $Y_\ell \in \{0, 1\}$  for  $\ell \in \mathcal{O} \cup \mathcal{E}$  sampled as indep. Bernoullis with

$$\Pr(Y_\ell(0) = 1 \mid \mathbf{x}_\ell) = \frac{1}{1 + e^{-\alpha - \beta^\top \mathbf{x}_\ell + \varepsilon_\ell}}, \quad \text{for } \beta = (1, 1, 1, 1, 1)^\top$$

for covariates  $X_\ell \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \mathbf{I}_5)$ ,  $\alpha$  chosen s.t. mean is 10%.

- Treatment variables  $W_j$  for  $j \in \mathcal{O}$  sampled via

$$\Pr(W_j = 1 \mid \mathbf{x}_j) = \frac{1}{1 + e^{-\gamma^\top \mathbf{x}_j}}, \quad \text{for } \gamma = (\sqrt{2}, \sqrt{2}, \sqrt{2}, 0, 0)^\top.$$

# Simulations Set-Up (II)

- Treatment effects
  - Define  $k = 1, \dots, 12$  strata based on first + second covariate
  - Assign  $\tau_k$ , stratum CATEs, via 3 treatment effect models:

$$\tau_k = T, \quad \tau_k = -T \times \frac{k}{K}, \quad \text{and} \quad \tau_k = T \times \left(\frac{k}{K}\right)^2$$

- $T$  chosen so that Cohen's D in ODB equals 0.5
- Simulation structure
  - Sample ODB data a single time. Correct via SIPW.
  - Compute RCT designs under different heuristics
  - Resample RCT units 5,000 times. For each sample, compute  $L_2$  error in estimating  $\tau$  using  $\hat{\tau}_r$ ,  $\kappa_2$ , and  $\kappa_{2+}$

# Simulations Set-Up (I)

- ODB has 20K units ( $j \in \mathcal{O}$ ). RCT has 1,000 ( $i \in \mathcal{E}$ )
- Untreated potential outcomes  $Y_\ell \in \{0, 1\}$  for  $\ell \in \mathcal{O} \cup \mathcal{E}$  sampled as indep. Bernoullis with

$$\Pr(Y_\ell(0) = 1 \mid \mathbf{x}_\ell) = \frac{1}{1 + e^{-\alpha - \beta^\top \mathbf{x}_\ell + \varepsilon_\ell}}, \quad \text{for } \beta = (1, 1, 1, 1, 1)^\top$$

for covariates  $X_\ell \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \mathbf{I}_5)$ ,  $\alpha$  chosen s.t. mean is 10%.

- Treatment variables  $W_j$  for  $j \in \mathcal{O}$  sampled via

$$\Pr(W_j = 1 \mid \mathbf{x}_j) = \frac{1}{1 + e^{-\gamma^\top \mathbf{x}_j}}, \quad \text{for } \gamma = (\sqrt{2}, \sqrt{2}, \sqrt{2}, 0, 0)^\top.$$

# Simulations Set-Up (II)

- Treatment effects
  - Define  $k = 1, \dots, 12$  strata based on first + second covariate
  - Assign  $\tau_k$ , stratum CATEs, via 3 treatment effect models:

$$\tau_k = T, \quad \tau_k = -T \times \frac{k}{K}, \quad \text{and} \quad \tau_k = T \times \left(\frac{k}{K}\right)^2$$

- $T$  chosen so that Cohen's D in ODB equals 0.5
- Simulation structure
  - Sample ODB data a single time. Correct via SIPW.
  - Compute RCT designs under different heuristics
  - Resample RCT units 5,000 times. For each sample, compute  $L_2$  error in estimating  $\tau$  using  $\hat{\tau}_r$ ,  $\kappa_2$ , and  $\kappa_{2+}$

# Idealized Case: All Covariates Measured

Est	Trt				Max Bias, $\Gamma$ Value				Oracle
		Eq.	Ney.	Naïve	1.0	1.1	1.2	1.5	
$\hat{\tau}_r$	c	100%	87%	91%	100%	96%	94%	94%	96%
$\kappa_2$		82%	48%	44%	52%	48%	47%	50%	42%
$\kappa_{2+}$		38%	28%	26%	26%	26%	26%	28%	23%
$\hat{\tau}_r$	l	100%	89%	92%	95%	94%	95%	97%	104%
$\kappa_2$		93%	66%	58%	58%	57%	60%	64%	50%
$\kappa_{2+}$		59%	51%	45%	43%	45%	47%	49%	33%
$\hat{\tau}_r$	q	100%	86%	91%	95%	98%	94%	92%	91%
$\kappa_2$		81%	47%	45%	52%	52%	50%	48%	41%
$\kappa_{2+}$		37%	29%	27%	28%	28%	30%	29%	25%

**Table 4:** Risk over 5,000 iterations of  $\hat{\tau}_r$ ,  $\kappa_2$ , and  $\kappa_{2+}$  in the case of no unmeasured confounding in the observational study. Risks are expressed as a percentage of the risk of  $\hat{\tau}_r$  using an equally allocated experiment, for each of the three treatment effect models.

# Simulations: Third Covariate Unmeasured

Est	Trt	Eq.	Ney.	Naïve	Worst Case, $\Gamma$ Value				Oracle
					1.0	1.1	1.2	1.5	
$\hat{\tau}_r$	c	100%	90%	90%	90%	92%	93%	95%	102%
$\kappa_2$		102%	81%	74%	72%	72%	72%	77%	69%
$\kappa_{2+}$		96%	80%	74%	71%	72%	72%	76%	67%
$\hat{\tau}_r$	$\ell$	100%	93%	93%	94%	95%	96%	96%	104%
$\kappa_2$		102%	85%	77%	75%	76%	77%	79%	73%
$\kappa_{2+}$		98%	84%	77%	75%	76%	76%	79%	71%
$\hat{\tau}_r$	q	100%	89%	90%	93%	92%	91%	96%	96%
$\kappa_2$		101%	74%	69%	68%	68%	67%	73%	66%
$\kappa_{2+}$		88%	72%	67%	66%	66%	65%	71%	63%

**Table 5:** Risk over 5,000 iterations of  $\hat{\tau}_r$ ,  $\kappa_2$ , and  $\kappa_{2+}$  under various experimental designs, in the case of unmeasured confounding in the observational study via failure to measure the third covariate.



## Combining these data: Locating the Problem

## How do we combine evidence from an RCT and an ODB?

This problem relates to several areas of research:

- Meta-analysis (Mueller et al., 2018; Prevost et al., 2000; Thompson et al., 2011)
- Transportability/generalizability (Stuart et al., 2011; Hartman et al., 2015; Bareinboim and Pearl, 2016)
- Causal inference (Kallus et al., 2018; Ghassami et al., 2022; Mooij et al., 2016)