

Shrinkage Estimation for Causal Inference and Experimental Design

Evan Rosenman[†], Guillaume Basse, Mike Baiocchi,
Art Owen, and Luke Miratrix

[†]Harvard University Data Science Initiative

November 17, 2022

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”

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)

Our Approach

We consider two problems:

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

Work in a stratified setting.

Shared conditional avg. treatment effect assumed within stratum

- Subject matter knowledge
- Algorithm for HTE (Wager and Athey, 2018; Hill, 2011)

Outline

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

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 for both populations:

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

Assumptions and Non-Assumptions

- ① Under F_O ,

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

No unconfoundedness assumption for observational study.

- ② Under F_R ,

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

- 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:

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

- Under mild conditions, we have

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

for bias $\boldsymbol{\xi}$ and diagonal covariance matrices Σ_r and Σ_o

- $\Sigma_r = \text{diag}(\sigma_{r1}^2, \dots, \sigma_{rK}^2)$ is estimable from the data
- $\boldsymbol{\xi}$ cannot be estimated, and estimates of Σ_o will be biased

Setup

- Collect our estimators into vectors:

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

- Under mild conditions, we have

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

for bias $\boldsymbol{\xi}$ and diagonal covariance matrices Σ_r and Σ_o

- $\Sigma_r = \text{diag}(\sigma_{r1}^2, \dots, \sigma_{rK}^2)$ is estimable from the data
- $\boldsymbol{\xi}$ cannot be estimated, and estimates of Σ_o will be biased
- Seek to design shrinkage estimator $\hat{\boldsymbol{\tau}} = f(\hat{\boldsymbol{\tau}}_r, \hat{\boldsymbol{\tau}}_o)$ to minimize expected L_2 loss (optionally weighted by \mathbf{W}),

$$\mathcal{L}(\hat{\boldsymbol{\tau}}, \boldsymbol{\tau}) = (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau})^T \mathbf{W} (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}).$$

Useful Prior Work

- **Shrinkage estimation:** a rich literature stretching back to multivariate normal mean estimation work of [Stein \(1956\)](#)
- [Green and Strawderman \(1991\)](#) and [Green et al. \(2005\)](#) propose estimators δ_1, δ_2 for shrinkage between ...
 - a normal, unbiased estimator (like $\hat{\tau}_r$), and
 - a biased estimator (like $\hat{\tau}_o$)
- **Key ideas**
 - Take convex combinations 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 $\hat{\tau}_o$
 - Estimators have bounded risk as $\hat{\tau}_o$ bias grows

Outline

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

Outline

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

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 λ 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 λ 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 λ 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 λ ,

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

A Note on λ_1^{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

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

λ_1^{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 λ_1^{URE} (I)

1 Define

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

κ_1 admits a testable condition under which it is guaranteed to reduce risk relative to $\hat{\tau}_r$.

Lemma (κ_1 Risk Guarantee)

Suppose $4 \max_k w_k \sigma_{rk}^2 < \sum_k w_k \sigma_{rk}^2$. Then κ_1 has risk strictly less than that of $\hat{\tau}_r$.

- Requires a dimension of at least $K = 4$.
- May require substantially larger K if high heteroscedasticity or non-uniform weights.

Useful Properties of λ_1^{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 λ_1^{URE} (III)

Theorem (κ_{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 $K \rightarrow \infty$, κ_{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 σ_{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,

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

Simulated Data Visualization

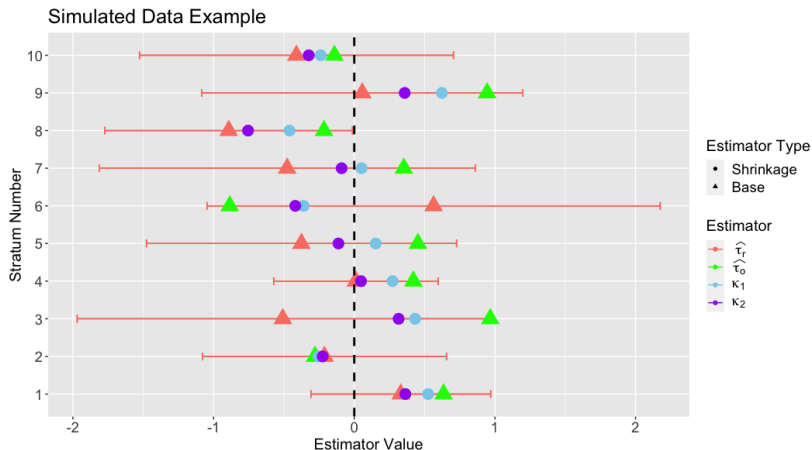


Figure 1: Simulated shrinkage between $\hat{\tau}_r$ and $\hat{\tau}_o$ with ten strata. 90% conf. sets for $\hat{\tau}_r$ in red, with κ_{1+} and κ_{2+} shown in circles.

Practical Considerations

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

Practical Considerations

- **Variance estimation:** In practice, Σ_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.

Practical Considerations

- **Variance estimation:** In practice, Σ_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.
- **Sensitivity analysis**
 - Marginal sensitivity model of [Tan \(2006\)](#) summarizes degree of unmeasured confounding by a single value, $\Gamma \geq 1$
 - Can “reverse engineer” implied confounding value Γ_{imp} when using a shrinker, via work of [Zhao et al. \(2019\)](#)
 - Evaluate Γ_{imp} to obtain a ✓ or ✗ for using shrinker

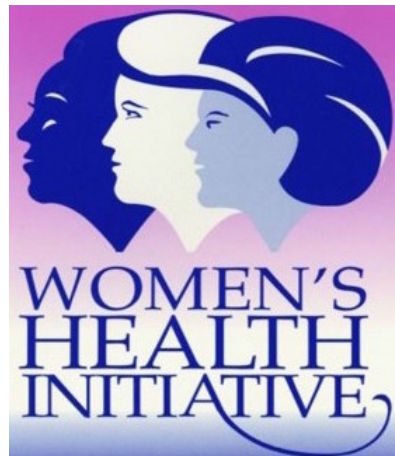
Outline

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

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 effects using entire RCT (16k units)
- Repeat 500 times:
 - Draw bootstrap samples:
 - 1,000 RCT units
 - Observational sample (50k units)
 - Compute L_2 loss for $\hat{\tau}_r, \kappa_{1+}, \kappa_{2+}, \delta_1, \delta_2$.
- Average loss over draws

Choice of Stratification Variables (I)

Stratify on two variables from WHI protocol:
age + history of cardiovascular disease ([Roehm, 2015](#)).

Age	Observational Study	RCT	RCT "Silver" Dataset
50-59	17,447 (33.0%)	5,491 (33.2%)	2,806 (33.9%)
60-69	23,030 (43.6%)	7,473 (45.2%)	3,689 (44.6%)
70-79	12,388 (23.4%)	3,573 (21.2%)	1,774 (21.5%)

CVD	Observational Study	RCT	RCT "Silver" Dataset
Yes	8,709 (16.5%)	1,828 (11.1%)	900 (10.9%)
No	44,156 (83.5%)	14,709 (88.9%)	7,369 (89.1%)

Choice of Stratification Variables (II)

Also include a variable unassociated with treatment effect:
solar irradiance (“Langley scatter”).

Exhibits no relevant correlations in the RCT gold dataset.

Langley Scatter	Observational Study	RCT	RCT “Silver” Dataset
300-325	15,599 (29.5%)	4,854 (29.4%)	2,411 (29.2%)
350	12,521 (23.7%)	3,917 (23.7%)	1,935 (23.4%)
375-380	5,841 (11.0%)	1,858 (11.2%)	934 (11.3%)
400-430	8,216 (15.5%)	2,585 (15.6%)	1,310 (15.8%)
475-500	10,688 (20.2%)	3,323 (20.1%)	1,679 (20.3%)

Results

Subgroup Variable(s)	# of Strata	Loss as % of $\hat{\tau}_r$ Loss			
		κ_{1+}	κ_{2+}	δ_1	δ_2
CVD	2	37.6%	36.9%	100.0%	100.0%
Age	3	37.3%	30.1%	61.5%	72.8%
Langley	5	29.4%	23.5%	40.0%	52.2%
CVD, Age	6	38.0%	38.2%	38.3%	82.4%
CVD, Langley	10	30.6%	32.5%	30.0%	87.2%
Age, Langley	15	22.4%	23.0%	22.5%	43.1%
Age, CVD, Langley	30	50.3%	50.3%	50.3%	78.4%

Outline

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

Outline

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

A New Setting: Design

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

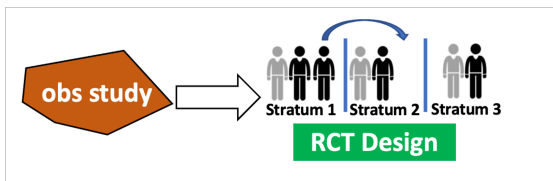
- Observational study already completed, $\hat{\tau}_o$ obtained.
- Designing a prospective RCT of n_r units
- Two objectives in mind
 - Aim to use a shrinker to combine $\hat{\tau}_r$ with $\hat{\tau}_o$. Design experiment to better complement ODB
 - Impose “guardrails” such that $\hat{\tau}_r$ still usable on its own

Defining the RCT “Design”

- **Goal:** choose an RCT allocation of treated and control counts per stratum, $\mathbf{d} = \{(n_{rkt}, n_{rk})\}_{k=1}^K$, such that

$$\sum_k n_{rkt} + n_{rk} = n_r$$

- Implies how to *recruit* (to get $n_{rk} = n_{rkt} + n_{rk}$ units for stratum k), ...
- and *assign* treatment via a simple random sample of n_{rkt} treated units from n_{rk} total units



Choice of Estimator: κ_2

We proceed with our estimator κ_{2+} from the prior section:

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

Consider $\hat{\tau}_o$ fixed. Compute exact risk of the estimator over randomness in $\hat{\tau}_r$:

$$\mathcal{R}(\kappa_2) = \frac{1}{K} \left(\text{Tr}(\Sigma_r) + \text{Tr}(\Sigma_r^2) \times \mathbb{E}_r \left(\frac{4(\hat{\tau}_r - \hat{\tau}_o)^\top \Sigma_r^4 (\hat{\tau}_r - \hat{\tau}_o)}{((\hat{\tau}_r - \hat{\tau}_o)^\top \Sigma_r^2 (\hat{\tau}_r - \hat{\tau}_o))^2} - \frac{\text{Tr}(\Sigma_r^2)}{(\hat{\tau}_r - \hat{\tau}_o)^\top \Sigma_r^2 (\hat{\tau}_r - \hat{\tau}_o)} \right) \right)$$

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 κ_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 $\mathcal{R}_2(\mathbf{d}, \mathbf{V}, \boldsymbol{\xi})$ as risk of κ_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 $\boldsymbol{\xi}$.

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 $\boldsymbol{\xi}$.

Estimating V : Updated Assumptions

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

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

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

Outline

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

1. Neyman Allocation

Under stronger Assumption 3, 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.

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

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 Γ -Level Unmeasured Confounding

- Can take a more pessimistic approach again using marginal sensitivity model of [Tan \(2006\)](#)
- Recall: 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 Γ -Level Unmeasured Confounding

- Can take a more pessimistic approach again using marginal sensitivity model of [Tan \(2006\)](#)
- Recall: 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

Guardrails

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

Recommend simple guardrails for designs:

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

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

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

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

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

Outline

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

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ || ≡ ↺ 🔍 ↻

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, λ_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 available at [arXiv:2002.06708](https://arxiv.org/abs/2002.06708)

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

References (I)

- Bao, Y. and Kan, R. (2013). On the moments of ratios of quadratic forms in normal random variables. *Journal of Multivariate Analysis*, 117:229–245.
- Bareinboim, E. and Pearl, J. (2016). Causal inference and the data-fusion problem. *Proceedings of the National Academy of Sciences*, 113(27):7345–7352.
- Ghassami, A., Shpitser, I., and Tchetgen, E. T. (2022). Combining experimental and observational data for identification of long-term causal effects. *arXiv preprint arXiv:2201.10743*.
- Green, E. J. and Strawderman, W. E. (1991). A James-Stein type estimator for combining unbiased and possibly biased estimators. *Journal of the American Statistical Association*, 86(416):1001–1006.
- Green, E. J., Strawderman, W. E., Amateis, R. L., and Reams, G. A. (2005). Improved estimation for multiple means with heterogeneous variances. *Forest Science*, 51(1):1–6.
- Hartman, E., Grieve, R., Ramsahai, R., and Sekhon, J. S. (2015). From sate to patt: combining experimental with observational studies to estimate population treatment effects. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 10:1111.
- Hill, J. L. (2011). Bayesian nonparametric modeling for causal inference. *Journal of Computational and Graphical Statistics*, 20(1):217–240.
- Kallus, N., Puli, A. M., and Shalit, U. (2018). Removing hidden confounding by experimental grounding. In *Advances in Neural Information Processing Systems*, pages 10888–10897.
- Li, K.-C. et al. (1985). From Stein's unbiased risk estimates to the method of generalized cross validation. *The Annals of Statistics*, 13(4):1352–1377.
- Mooij, J. M., Magliacane, S., and Claassen, T. (2016). Joint causal inference from multiple contexts. *arXiv preprint arXiv:1611.10351*.
- Mueller, M., D'Addario, M., Egger, M., Cevallos, M., Dekkers, O., Mugglin, C., and Scott, P. (2018). Methods to systematically review and meta-analyse observational studies: a systematic scoping review of recommendations. *BMC Medical Research Methodology*, 18(1):44.
- Prevost, T. C., Abrams, K. R., and Jones, D. R. (2000). Hierarchical models in generalized synthesis of evidence: an example based on studies of breast cancer screening. *Statistics in Medicine*, 19(24):3359–3376.

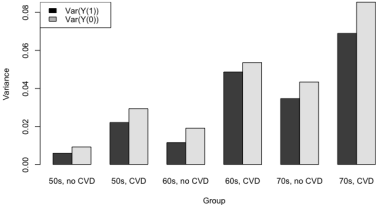
References (II)

- Roehm, E. (2015). A reappraisal of Women's Health Initiative estrogen-alone trial: long-term outcomes in women 50–59 years of age. *Obstetrics and Gynecology International*, 2015.
- Rosenbaum, P. R. (1987). Sensitivity analysis for certain permutation inferences in matched observational studies. *Biometrika*, 74(1):13–26.
- Stein, C. (1956). Inadmissibility of the usual estimator for the mean of a multivariate normal distribution. Technical report, Stanford University Stanford United States.
- Stein, C. M. (1981). Estimation of the mean of a multivariate normal distribution. *The Annals of Statistics*, pages 1135–1151.
- Stuart, E. A., Cole, S. R., Bradshaw, C. P., and Leaf, P. J. (2011). The use of propensity scores to assess the generalizability of results from randomized trials. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 174(2):369–386.
- Tan, Z. (2006). A distributional approach for causal inference using propensity scores. *Journal of the American Statistical Association*, 101(476):1619–1637.
- Thompson, S., Ekelund, U., Jebb, S., Lindroos, A. K., Mander, A., Sharp, S., Turner, R., and Wilks, D. (2011). A proposed method of bias adjustment for meta-analyses of published observational studies. *International journal of epidemiology*, 40(3):765–777.
- Wager, S. and Athey, S. (2018). Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113(523):1228–1242.
- Xie, X., Kou, S., and Brown, L. D. (2012). Sure estimates for a heteroscedastic hierarchical model. *Journal of the American Statistical Association*, 107(500):1465–1479.
- Zhao, Q., Small, D. S., and Bhattacharya, B. B. (2019). Sensitivity analysis for inverse probability weighting estimators via the percentile bootstrap. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*.

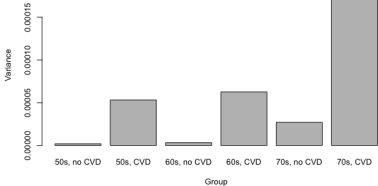
Appendices

Sample Designs

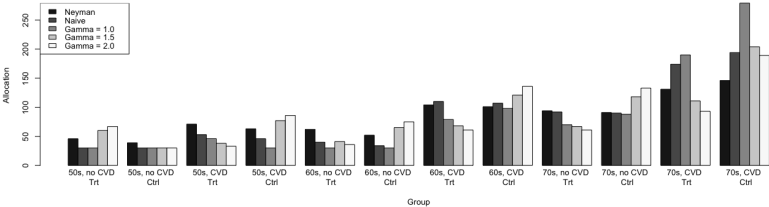
Potential Outcome Variance by Stratum



Observational Estimator Variance by Stratum



Allocations by Stratum Under Different Schemes



- 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
Age	60.78	64.72	-0.56	63.06	63.33	-0.04
BMI	25.55	27.11	-0.25	26.71	26.62	0.00
Physical functioning	85.23	79.58	0.26	81.15	81.23	0.03
Age at menopause	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
Before Strat.	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%	
After Strat.	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.

		Never Smoked	Past Smoker	Current Smoker	SD
Before Stratifying	Treated	48.7%	46.2%	5.1%	0.11
	Control	52.3%	41.1%	6.6%	
After Stratifying	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 ξ 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.

δ_1 is designed for precision-weighted loss, but outperforms δ_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 κ_{1+}

- Recall: λ_1^{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 κ_{1+} !
- Key idea:** frame in context of sensitivity model of [Tan \(2006\)](#)

- Marginal sensitivity model of Tan (2006) summarizes degree of unmeasured confounding by a single value, $\Gamma \geq 1$
 - Γ 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 Γ
 - Implicitly maps Γ to a worst-case bias $\xi(\Gamma)$ and variance $\Sigma_O(\Gamma)$
 - Under some assumptions, allows us to obtain worst-case estimate of λ_{opt} as a function of Γ , which we call $\lambda(\Gamma)$

Relating the Models

- **Intuition:** larger Γ (confounding parameter) \implies optimal weight λ_{opt} is smaller
- Let $\Gamma_{\text{imp}} = \sup\{\Gamma : \lambda(\Gamma) > \lambda_1^{\text{URE}}\}$
 - Largest value Γ for which the optimal shrinkage factor $\lambda(\Gamma)$ is greater than our shrinkage parameter λ_1^{URE} .
- Γ_{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 $X_\ell \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \mathbf{I}_5)$, α 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 τ_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 τ using $\hat{\tau}_r$, κ_2 , and κ_{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 $X_\ell \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \mathbf{I}_5)$, α 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 τ_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 τ using $\hat{\tau}_r$, κ_2 , and κ_{2+}

Idealized Case: All Covariates Measured

Est	Trt				Max Bias, Γ 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%
κ_2		82%	48%	44%	52%	48%	47%	50%	42%
κ_{2+}		38%	28%	26%	26%	26%	26%	28%	23%
$\hat{\tau}_r$	l	100%	89%	92%	95%	94%	95%	97%	104%
κ_2		93%	66%	58%	58%	57%	60%	64%	50%
κ_{2+}		59%	51%	45%	43%	45%	47%	49%	33%
$\hat{\tau}_r$	q	100%	86%	91%	95%	98%	94%	92%	91%
κ_2		81%	47%	45%	52%	52%	50%	48%	41%
κ_{2+}		37%	29%	27%	28%	28%	30%	29%	25%

Table 4: Risk over 5,000 iterations of $\hat{\tau}_r$, κ_2 , and κ_{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, Γ Value				Oracle
					1.0	1.1	1.2	1.5	
$\hat{\tau}_r$	c	100%	90%	90%	90%	92%	93%	95%	102%
κ_2		102%	81%	74%	72%	72%	72%	77%	69%
κ_{2+}		96%	80%	74%	71%	72%	72%	76%	67%
$\hat{\tau}_r$	ℓ	100%	93%	93%	94%	95%	96%	96%	104%
κ_2		102%	85%	77%	75%	76%	77%	79%	73%
κ_{2+}		98%	84%	77%	75%	76%	76%	79%	71%
$\hat{\tau}_r$	q	100%	89%	90%	93%	92%	91%	96%	96%
κ_2		101%	74%	69%	68%	68%	67%	73%	66%
κ_{2+}		88%	72%	67%	66%	66%	65%	71%	63%

Table 5: Risk over 5,000 iterations of $\hat{\tau}_r$, κ_2 , and κ_{2+} under various experimental designs, in the case of unmeasured confounding in the observational study via failure to measure the third covariate.