

# To Link or Not to Link? Estimating Long-run Treatment Effects from Historical Data

Francis J. DiTraglia<sup>1</sup>   Camilo García-Jimeno<sup>2</sup>   Ezra Karger<sup>2</sup>

<sup>1</sup>Department of Economics, University of Oxford

<sup>2</sup>Federal Reserve Bank of Chicago

December 14th, 2024

The views expressed in this talk are those of the authors and do not necessarily reflect the position of the Federal Reserve Bank of Chicago or the Federal Reserve System.

# The Linking Problem

Aizer et al (2016; AER)

- ▶ Long-term effects of “Mothers’ Pension Program” on adult outcomes
- ▶ Outcomes from 1940 U.S. Census; Treatment from 1911-1935 Admin. Records

Abramitzky, Boustan & Eriksson (ABE) Algorithm

- ▶ Form linked dataset of treatments/outcomes for a subset of individuals
- ▶ Exact or near agreement of linking variables: name, sex, race, age, state of birth
- ▶ Abramitzky et al. (2021; JEL) “Automated Linking of Historical Data”

# Methodological Challenges

## Low Match Rates

- ▶ In typical applications of ABE only  $\approx 20\%$  of observations are matched
- ▶ Women and minorities typically excluded altogether: harder to match
- ▶ Inefficient estimation; potentially unrepresentative sample

## Noisy Linking Variables

- ▶ Variables typically used for linking are known to be measured with error
- ▶ Implicitly acknowledged in ABE algorithm: permits “near” matches
- ▶ Ignored in practice: linked dataset analyzed as though it has no spurious matches

Learn  $\beta$  from  $Y_i = X_i'\beta + U$  where  $\mathbb{E}(U_i|X_i) = 0$

	index	$W$	$Y$		index	$W$	$X$	$Y$
DF <sub>y</sub>	1	$\widetilde{W}_1$	$\widetilde{Y}_1$	DF <sub>x</sub>	1	$W_1$	$X_1$	$Y_1$
	$\vdots$	$\vdots$	$\vdots$		$\vdots$	$\vdots$	$\vdots$	$\vdots$
	$j$	$\widetilde{W}_j$	$\widetilde{Y}_j$		$j'$	$W_{j'}$	$X_{j'}$	$Y_{j'}$
	$\vdots$	$\vdots$	$\vdots$		$\vdots$	$\vdots$	$\vdots$	$\vdots$
	$n$	$\widetilde{W}_n$	$\widetilde{Y}_n$		$n$	$W_n$	$X_n$	$Y_n$

- Blue means observed; Red means unobserved
- $W$  observed jointly with outcomes in DF<sub>y</sub> and jointly with treatments in DF<sub>x</sub>

## The Linking Matrix

$\mathbf{L} = [\ell_{jj'}]$  is an  $(n \times n)$  permutation matrix where

$$\ell_{jj'} = \begin{cases} 1 & \text{if record } j \text{ in } \text{DF}_y \text{ matches with record } j' \text{ in } \text{DF}_x \\ 0 & \text{otherwise.} \end{cases}$$

Since the  $j^{\text{th}}$  row of  $\mathbf{L}$  is  $[\ell_{j1} \ \ell_{j2} \ \cdots \ \ell_{jN}]$ , it follows that  $\tilde{\mathbf{Y}} = \mathbf{L}\mathbf{Y}$

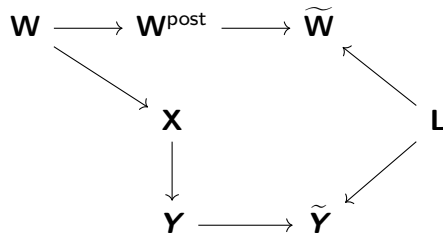
$$\tilde{Y}_j = \begin{pmatrix} Y_1 & \text{if } \ell_{j1} = 1 \\ Y_2 & \text{if } \ell_{j2} = 1 \\ \vdots & \\ Y_N & \text{if } \ell_{jN} = 1 \end{pmatrix} = \sum_{j'=1}^N \ell_{jj'} Y_{j'}.$$

# Fundamental Decomposition

(A)  $Y \perp\!\!\!\perp L \mid (X, W, \widetilde{W})$

(B)  $Y \perp\!\!\!\perp (W, \widetilde{W}) \mid X$

(C)  $L \perp\!\!\!\perp X \mid (W, \widetilde{W})$



$$\begin{aligned}
 \mathbb{E}(\widetilde{Y} | X, W, \widetilde{W}) &= \mathbb{E}_{L|X, W, \widetilde{W}} \left[ \mathbb{E}(LY | X, W, \widetilde{W}, L) \right] \\
 &= \mathbb{E}_{L|X, W, \widetilde{W}} \left[ L \mathbb{E}(Y | X, W, \widetilde{W}) \right] \\
 &= \mathbb{E}(L | X, W, \widetilde{W}) \mathbb{E}(Y | X, W, \widetilde{W}) \\
 &= Q(W, \widetilde{W})(X\beta)
 \end{aligned}$$

# Estimators that Regress $\widetilde{\mathbf{Y}}$ on $\mathbf{QX} \equiv$ Imputed $\widetilde{\mathbf{X}}$

## Unique Match (UM)

Run ABE; fill  $\mathbf{Q}$  with ones and zeros to indicate unique matches; drop many  $\widetilde{Y}_j$

## Poirier & Ziebarth (PZ)

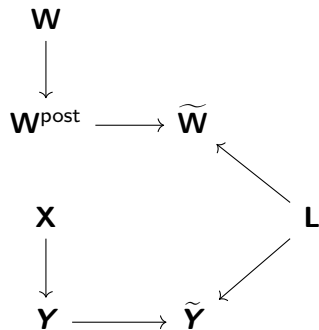
Run ABE; form groups of “potential matches”; give equal weight to each in  $\mathbf{Q}$

## Probabilistic Linking (PL)

- ▶ Compute  $\mathbf{Q}(\mathbf{W}, \widetilde{\mathbf{W}}) \equiv \mathbb{E}(\mathbf{L}|\mathbf{W}, \widetilde{\mathbf{W}})$  using Bayes' Theorem
- ▶ Tempting simplification:  $\mathbb{P}(\ell_{jj'} = 1 | \mathbb{1}\{\widetilde{W}_j = W_{j'}\})$

# Simulation Design

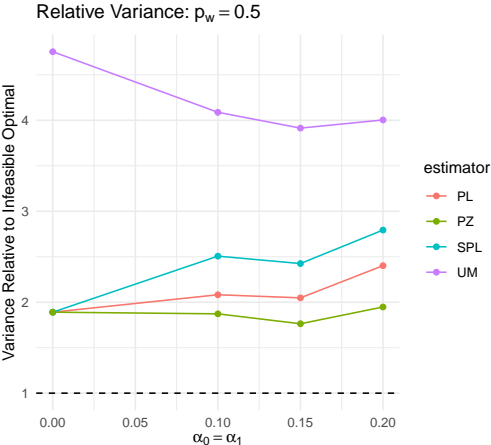
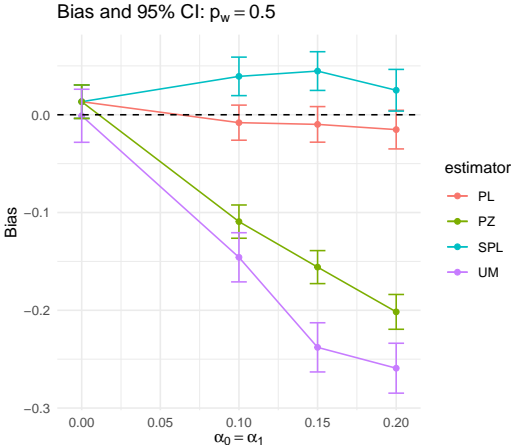
- (i)  $X_i \sim \text{iid Uniform}(0, 1)$
- (ii)  $Y_i|X_i \sim \text{iid Normal}(X_i, \sigma^2 = 1)$
- (iii)  $W_i \sim \text{iid Bernoulli}(p_w)$
- (iv)  $W_i^{\text{post}}|W_i \sim \text{iid Bernoulli}((1 - W_i)\alpha_0 + W_i(1 - \alpha_1))$



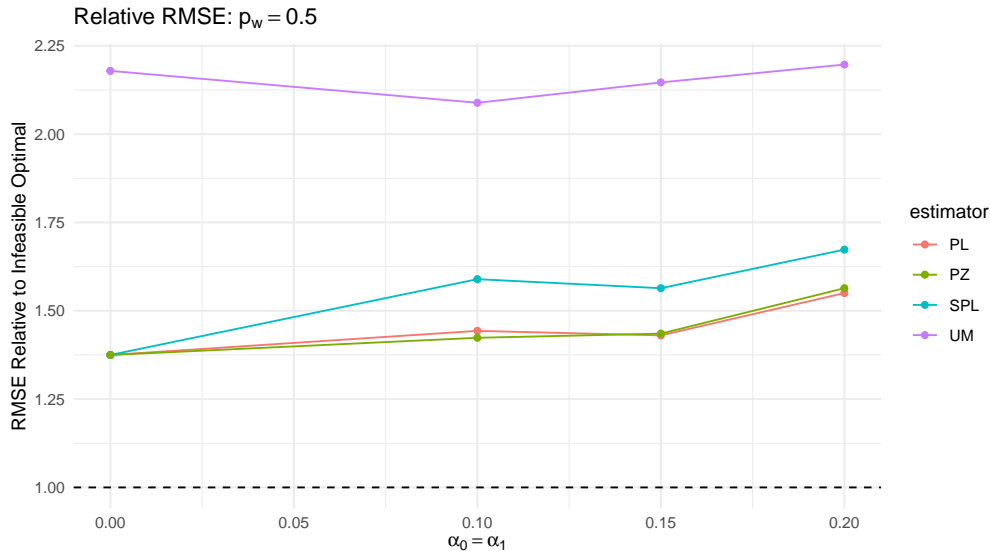
- ▶ Mis-classification probabilities  $\alpha_0$  and  $\alpha_1$
- ▶  $(\tilde{W}, \tilde{Y})$  random re-ordering of the rows of  $(W^{\text{post}}, Y)$  within “blocks”
- ▶ 50 blocks (“states”) each containing  $[2 + \text{Poisson}(1)]$  individuals ( $n \approx 150$ )



# Simulation Results



# Simulation Results



## Next Steps

- ▶ Head-to-head comparisons using real data: how different are the results?
- ▶ Our “full” PL estimator works well, but is hard to scale up. Approximations?
- ▶ Connection with TS2SLS using mis-classified instruments.
- ▶ Additional complication: “treatment” is often a function of  $X$  and  $W$
- ▶ Start thinking more carefully about sample selection / heterogeneity