

# Distributional Treatment Effect with Latent Rank Invariance

Myungkou Shin

University of Surrey

McGill University Job Market Seminar

January 9, 2026

## Distributional treatment effect

Motivation: in discussions of treatment effect heterogeneity,  
*distribution* of treatment effect lies at the core.

Goal: estimate distributional treatment effect parameters  
under intuitive and easily interpretable assumptions.

## Distributional treatment effect

Motivation: in discussions of treatment effect heterogeneity, *distribution* of treatment effect lies at the core.

Goal: estimate distributional treatment effect parameters under intuitive and easily interpretable assumptions.

A potential outcome setup with a binary treatment: with  $D_i \in \{0, 1\}$ ,

$$Y_i = D_i \cdot Y_i(1) + (1 - D_i) \cdot Y_i(0)$$

We only observe  $Y_i(1)$  or  $Y_i(0)$  for a given unit  $i$ .

## Distributional treatment effect

Motivation: in discussions of treatment effect heterogeneity,  
*distribution* of treatment effect lies at the core.

Goal: estimate distributional treatment effect parameters  
under intuitive and easily interpretable assumptions.

A potential outcome setup with a binary treatment: with  $D_i \in \{0, 1\}$ ,

$$Y_i = D_i \cdot Y_i(1) + (1 - D_i) \cdot Y_i(0)$$

We only observe  $Y_i(1)$  or  $Y_i(0)$  for a given unit  $i$ .

Existing frameworks mostly focus on some summary measures of treatment effect:

e.g.,	Average treatment effect on treated units (ATT)	$\mathbf{E} [Y_i(1) - Y_i(0)   D_i = 1]$
	Average treatment effect (ATE)	$\mathbf{E} [Y_i(1) - Y_i(0)]$
	Quantile treatment effect (QTE( $\tau$ ))	$F_{Y(1)}^{-1}(\tau) - F_{Y(0)}^{-1}(\tau)$

## Distributional treatment effect

Summary measures are invariant the dependence structure between  $Y_i(1)$  and  $Y_i(0)$ .

Can't answer questions on treatment effect distribution or joint distribution of  $Y_i(1)$  and  $Y_i(0)$ :

## Distributional treatment effect

Summary measures are invariant the dependence structure between  $Y_i(1)$  and  $Y_i(0)$ .

Can't answer questions on treatment effect distribution or joint distribution of  $Y_i(1)$  and  $Y_i(0)$ :

1. Fairness concern: Levy and Markowitz (1979); Epstein and Segal (1992) and more

“Aggregate treatment effect with inequality aversion.”

“I'd like to look at treatment effect heterogeneity w.r.t. baseline.”

- Variance-adjusted ATE  $\mathbf{E}[Y_i(1) - Y_i(0)] + \frac{\gamma}{2} \text{Var}(Y_i(1) - Y_i(0)).$
- Conditional ATE given  $Y_i(0) \in \mathcal{Y}$   $\mathbf{E}[Y_i(1) - Y_i(0) | Y_i(0) \in \mathcal{Y}].$

## Distributional treatment effect

Summary measures are invariant the dependence structure between  $Y_i(1)$  and  $Y_i(0)$ .

Can't answer questions on treatment effect distribution or joint distribution of  $Y_i(1)$  and  $Y_i(0)$ :

1. Fairness concern: Levy and Markowitz (1979); Epstein and Segal (1992) and more

“Aggregate treatment effect with inequality aversion.”

“I'd like to look at treatment effect heterogeneity w.r.t. baseline.”

2. Probabilistic guarantee: Reeve et al. (2023); Chernozhukov et al. (2025)

“How many people are better off under treatment?”

“What is  $100 \cdot \alpha\%$  worst-case scenario?”

- Share of winners  $1 - F_{Y(1)-Y(0)}(0).$

-  $\alpha$ -th quantile of treatment effect  $F_{Y(1)-Y(0)}^{-1}(\alpha).$

## Distributional treatment effect

Summary measures are invariant the dependence structure between  $Y_i(1)$  and  $Y_i(0)$ .

Can't answer questions on treatment effect distribution or joint distribution of  $Y_i(1)$  and  $Y_i(0)$ :

1. Fairness concern: Levy and Markowitz (1979); Epstein and Segal (1992) and more

“Aggregate treatment effect with inequality aversion.”

“I'd like to look at treatment effect heterogeneity w.r.t. baseline.”

2. Probabilistic guarantee: Reeve et al. (2023); Chernozhukov et al. (2025)

“How many people are better off under treatment?”

“What is  $100 \cdot \alpha\%$  worst-case scenario?”

3. Voluntary take-up: Heckman and Vytlacil (2005); Mogstad et al. (2018) and more

“How many people would opt into treatment at the cost of  $c$ .”

“Policy-relevant treatment effect when people can opt out with full information.”

- Take-up at cost  $c$   $F_{Y(1)-Y(0)}(c)$ .

- Conditional ATE given  $Y_i(1) \geq Y_i(0)$   $\mathbf{E}[Y_i(1) - Y_i(0)|Y_i(1) \geq Y_i(0)]$ .



## Distributional treatment effect

Existing approaches on  $F_{Y(1)-Y(0)}$

- Partial identification: put bounds on  $F_{Y(1)-Y(0)}(\delta)$ .

Heckman et al. (1997); Fan and Park (2010); Firpo and Ridder (2019); Frandsen and Lefgren (2021); Kaji and Cao (2023) and more

Makarov bound; optimal transport with additional constraints

⇒ Bounds often uninformatively large.

- Independence: assume  $Y_i(1) \perp\!\!\!\perp Y_i(0)$  or  $Y_i(0) \perp\!\!\!\perp (Y_i(1) - Y_i(0))$

Heckman et al. (1997); Carneiro et al. (2003); Wu and Perloff (2006); Noh (2023)

Multiplication and integration; deconvolution.

⇒ Rely on restrictive independence or functional form assumptions.

## Distributional treatment effect

Assume a **conditional independence** framework: Carneiro et al. (2003)

- Nonparametric identification and estimation.

Point identify and estimate **distributional treatment effect (DTE)**  $\theta$  such that

$$\mathbf{E} [m(Y_i(1), Y_i(0); \theta)] = 0.$$

- Moment-identified DTE parameter  $\theta$  include:

$\text{Var}(Y_i(1) - Y_i(0))$ ,  $F_{Y(1)-Y(0)}(\delta)$  for some  $\delta$ ,  $F_{Y(0), Y(1)}(y, y')$  for some  $(y, y')$  and many more.

## Conditional independence

In this paper, I assume a latent variable  $U_i$  such that

$$Y_i(1) \perp\!\!\!\perp Y_i(0) \mid U_i.$$

$U_i$  models individual-level heterogeneity and explains the dependence between  $Y_i(1)$  and  $Y_i(0)$ .

$$\begin{aligned}\Pr \{Y_i(0) \leq y, Y_i(1) \leq y'\} &= \mathbf{E} [\Pr \{Y_i(0) \leq y, Y_i(1) \leq y' | U_i\}] \\ &= \mathbf{E} [\Pr \{Y_i(0) \leq y | U_i\} \cdot \Pr \{Y_i(1) \leq y' | U_i\}].\end{aligned}$$

## Conditional independence

In this paper, I assume a latent variable  $U_i$  such that

$$Y_i(1) \perp\!\!\!\perp Y_i(0) \mid U_i.$$

$U_i$  models individual-level heterogeneity and explains the dependence between  $Y_i(1)$  and  $Y_i(0)$ .

$$\begin{aligned}\Pr \{Y_i(0) \leq y, Y_i(1) \leq y'\} &= \mathbf{E} [\Pr \{Y_i(0) \leq y, Y_i(1) \leq y' | U_i\}] \\ &= \mathbf{E} [\Pr \{Y_i(0) \leq y | U_i\} \cdot \Pr \{Y_i(1) \leq y' | U_i\}].\end{aligned}$$

Once I identify

1. conditional dist. of  $Y_i(1)$  given  $U_i$
2. conditional dist. of  $Y_i(0)$  given  $U_i$
3. marginal dist. of  $U_i$

I identify  $F_{Y(0), Y(1)}$  and thus any moment-identified DTE parameter  $\theta$ .

To identify the distributions, I assume two proxy variables  $X_i$  and  $Z_i$ , which shift  $U_i$ .

## Conditional independence: setup

An econometrician observes  $\{Y_i, D_i, X_i, Z_i\}_{i=1}^n$ :

$$Y_i = D_i \cdot Y_i(1) + (1 - D_i) \cdot Y_i(0).$$

There exists a latent variable  $U_i$ .

$Y_i(1), Y_i(0), X_i, Z_i, U_i \in \mathbb{R}$  and  $D_i \in \{0, 1\}$ .

$(Y_i(1), Y_i(0), D_i, X_i, Z_i, U_i) \sim iid.$

## Conditional independence: setup

An econometrician observes  $\{Y_i, D_i, X_i, Z_i\}_{i=1}^n$ :

$$Y_i = D_i \cdot Y_i(1) + (1 - D_i) \cdot Y_i(0).$$

There exists a latent variable  $U_i$ .

$Y_i(1), Y_i(0), X_i, Z_i, U_i \in \mathbb{R}$  and  $D_i \in \{0, 1\}$ .

$(Y_i(1), Y_i(0), D_i, X_i, Z_i, U_i) \sim iid.$

### **Assumptions 1 and 2.** (*conditional independence*)

$Y_i(1), Y_i(0), X_i, (D_i, Z_i)$  are mutually independent of each other given  $U_i$ .

- Only one proxy  $Z_i$  may depend on  $D_i$  given  $U_i$ .
- Can be connected to proximal inference and nonclassical measurement error literature. [more](#)
- $Y_i(1) \perp\!\!\!\perp Y_i(0) \mid U_i$  gives us key identifying power. [more](#)

Conditional independence: what is  $U_i$  and where do we find  $X_i$  and  $Z_i$ ?

### 1. Measurement error model

In some empirical contexts, natural interpretation on  $U_i$ .

For example,  $D_i$  is early childhood intervention and  $Y_i$  is cognitive development score.

Then, " $U_i =$  the innate ability of a child,

$(X_i, Z_i) =$  repeated measures of the innate ability, such as test scores"

Carneiro et al. (2003); Cunha et al. (2010); Attanasio et al. (2020) and more.

Conditional independence: what is  $U_i$  and where do we find  $X_i$  and  $Z_i$ ?

## 1. Measurement error model

In some empirical contexts, natural interpretation on  $U_i$ .

For example,  $D_i$  is early childhood intervention and  $Y_i$  is cognitive development score.

Then, “ $U_i$  = the innate ability of a child,

$(X_i, Z_i)$  = repeated measures of the innate ability, such as test scores”

Carneiro et al. (2003); Cunha et al. (2010); Attanasio et al. (2020) and more.

## 2. Hidden Markov model for panel data

A hidden Markov model for potential outcomes:

$$Y_{it}(d) = g_d(V_{it}, \varepsilon_{it}^d)$$

and  $\{V_{it}\}_t$  is first-order Markovian.

Then, “ $U_i$  = the contemporaneous common shock ( $V_{i2}$ ),

$(X_i, Z_i)$  = past and future outcomes ( $Y_{i1}$  and  $Y_{i3}$ ).”

Kasahara and Shimotsu (2009); Hu and Shum (2012); Deaner (2023) and more



### Identification

1. Conditional independence framework:  $Y_i(1), Y_i(0), X_i, (D_i, Z_i) \mid U_i \sim \text{ind.}$
2. Apply diagonalization (Hu and Schennach, 2008) to untreated and treated subpopulations.
3. Connect the two decomposition results to identify  $F_{Y(0), Y(1)}$ .

### Estimation

finite support for  $U_i \Rightarrow$  conditional independence becomes finite mixture

1. First-step: nonnegative matrix factorization (NMF) for finite mixture.  
a new estimator; improved finite sample performance.
2. Second-step: plug-in GMM for DTE.  
first-step NMF as nuisance parameters.  
asymptotic normality thanks to Neyman orthogonality.

### **Nonclassical measurement error/proximal inference/finite mixture**

- Estimation mostly relies on diagonalization

Hu (2008); Kasahara and Shimotsu (2009); Bonhomme et al. (2016) and more.

⇒ finite mixture estimator with additional regularization; better finite sample performance.  
 $\sqrt{n}$ -consistency and orthogonalization procedure proposed.

### **Distributional treatment effect**

- Mostly focus on partial identification.

Fan and Park (2010); Fan et al. (2014); Firpo and Ridder (2019); Frandsen and Lefgren (2021); Kaji and Cao (2023) and more.

- A few notable point identification exceptions:

Heckman et al. (1997); Carneiro et al. (2003); Wu and Perloff (2006); Noh (2023).

⇒ DTE estimator not relying on parametric distributions nor unconditional independence.  
Huge information gain compared to partial bounds.

## Preview of results

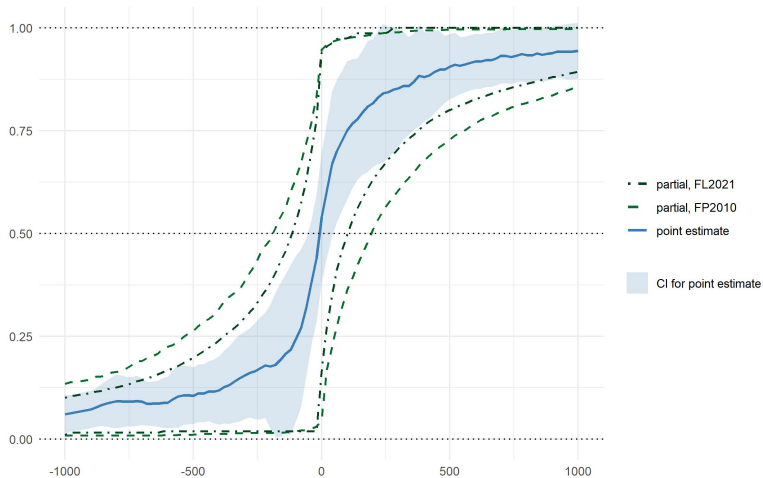


Figure 1: Marginal distribution of  $Y_i(1) - Y_i(0)$ .

## Identification

## Identification

The goal is to identify

$$\Pr \{Y_i(0) \leq y, Y_i(1) \leq y'\} = \mathbf{E} [\Pr \{Y_i(0) \leq y | U_i\} \cdot \Pr \{Y_i(1) \leq y' | U_i\}] .$$

Then, any DTE parameter  $\theta$  s.t.  $\mathbf{E} [m(Y_i(1), Y_i(0); \theta)] = 0$  is identified.

Identification strategy:

1. Identify the conditional distribution of  $Y_i(d) \mid U_i$ , within each subpopulation.
2. Identify the distribution of  $U_i$ .
3. Integrate out  $U_i$  from the conditional distribution of  $(Y_i(1), Y_i(0)) \mid U_i$ .

## Identification: diagonalization à la Hu and Schennach (2008)

Given  $(y, d)$ , construct

$$\mathbf{H}_d(y) = \begin{pmatrix} \Pr\{Y_i = y, X_i = x^1 | D_i = d, Z_i = z^1\} & \cdots & \Pr\{Y_i = y, X_i = x^1 | D_i = d, Z_i = z^J\} \\ \vdots & \ddots & \vdots \\ \Pr\{Y_i = y, X_i = x^{M_X} | D_i = d, Z_i = z^1\} & \cdots & \Pr\{Y_i = y, X_i = x^{M_X} | D_i = d, Z_i = z^J\} \end{pmatrix}.$$

$Y_i, X_i, Z_i$  can be discretized when continuous.

Suppose  $U_i$  is discrete. Under [Assumptions 1-2](#) that  $Y_i(1), Y_i(0), X_i, (D_i, Z_i) \mid U_i \sim \text{ind}$

$$\mathbf{H}_d(y) = \Gamma_X \cdot \Delta_d(y) \cdot \Lambda_d$$

where  $\Gamma_X = \left( \Pr\{X_i = x^m | U_i = u^k\} \right)_{m,k}$   
 $\Delta_d(y) = \text{diag}\left( \Pr\{Y_i(d) = y | U_i = u^k\} \right)_k$   
 $\Lambda_d = \left( \Pr\{U_i = u^k | D_i = d, Z_i = z^j\} \right)_{k,j}.$

## Identification: diagonalization à la Hu and Schennach (2008)

$$\sum_{y'} \mathbf{H}_d(y') = \Gamma_X \cdot \left( \sum_{y'} \Delta_d(y') \right) \cdot \Lambda_d = \Gamma_X \cdot \Lambda_d.$$

When  $\Gamma_X \cdot \Lambda_d$  is invertible, we get

$$\mathbf{H}_d(y) \left( \sum_{y'} \mathbf{H}_d(y') \right)^{-1} = \Gamma_X \cdot \Delta_d(y) \cdot \Lambda_d \cdot \left( \Gamma_X \cdot \Lambda_d \right)^{-1} = \Gamma_X \cdot \Delta_d(y) \cdot \left( \Gamma_X \right)^{-1}$$

Eigenvalue decomposition finds  $\Delta_d(y)$  and  $\Gamma_X$  up to sign and scale.

Repeating this across  $y$  completes identification: Hu (2008)

Hu and Schennach (2008) develops its counterpart for continuous  $U_i$ . [more](#)

Conditional densities  $f_{X|U}$ ,  $f_{Y(d)|U}$  are identified.

$d$  is fixed; needs to extend Hu and Schennach (2008) to potential outcome setup.

## Identification: sketchy of proof

1. Apply Hu and Schennach (2008) to the two subpopulations:

$$(\Gamma_X, \{\Delta_0(y)\}_y) \quad \text{and} \quad (\Gamma_X, \{\Delta_1(y)\}_y).$$

Labelings on  $U_i$  are connected using  $\Gamma_X$  since  $X_i \perp\!\!\!\perp D_i \mid U_i$ .

2.  $\{\Delta_0(y)\}_y$  give us  $f_{Y(0)|U}$  and  $\{\Delta_1(y)\}_y$  give us  $f_{Y(1)|U}$ .
3. Given  $\Gamma_X$ , we get

$$\Lambda_d = \left( f_{U|D,Z}(u|d,z) \right)_{u,z} = (\Gamma_X)^+ \sum_{y'} \mathbf{H}_d(y').$$

$\Lambda_0, \Lambda_1$  and the observed distribution of  $(D_i, Z_i)$  identifies the distribution of  $U_i$ :

$$f_U(u) = \mathbf{E} \left[ f_{U|D,Z}(u|D_i, Z_i) \right].$$

$\Rightarrow f_{Y(1)|U}, f_{Y(0)|U}, f_U$  are identified.  $Y_i(1) \perp\!\!\!\perp Y_i(0) \mid U_i$  identifies  $F_{Y(0), Y(1)}$ .



## Identification: identifying assumptions

**Assumption 3/4.** full rank/completeness of  $f_{X|Z}$  when  $U_i$  is discrete/continuous: A3 A4

“Both of the proxy variables are fully informative about the latent variable  $U_i$ .”

In the case of continuous  $U_i$ ,

**Assumption 5.**  $\mathbf{E}[Y_i(1)|U_i = u]$  or  $\mathbf{E}[Y_i(0)|U_i = u]$  is strictly increasing in  $u$ .

“For either  $d = 0$  or  $1$ , the latent variable  $U_i$  is the rank of  $\mathbf{E}[Y_i(d)|U_i]$ .”

Motivated from quantile treatment effect/IV literature:

Chernozhukov and Hansen (2005, 2006); Athey and Imbens (2006); Callaway and Li (2019) and more.

Rank invariance assumes a deterministic relationship between  $Y_i(1)$  and  $Y_i(0)$ .

In my framework, only the systemic parts of  $Y_i(1)$  and  $Y_i(0)$  are connected through  $U_i$ :  
*latent rank invariance*.

Assumption 5 says  $U_i$  is the latent rank for either  $Y_i(1)$  or  $Y_i(0)$ .

**Theorem 1.**

Assumptions 1-3 or Assumptions 1-2, 4-5 hold.

Then, the distribution of  $(Y_i(1), Y_i(0), D_i, X_i, Z_i)$  is identified.

Any moment-identified DTE parameter is trivially identified.

- Marginal distribution of treatment effect  $\theta = F_{Y(1)-Y(0)}(\delta)$  as a working example.

Can be extended to multidimensional  $U_i$ . [more](#)

- The dimension of  $X_i$  and  $Z_i$  needs to be at least equal to the dimension of  $U_i$ .

## Implementation

## Implementation: finite support assumption

I assume  $U_i \in \{u^1, \dots, u^K\}$  with  $K < \infty$ . choice of  $K$

Discretization as approximation: Bonhomme et al. (2022) and more.

Reasoning behind the finite support assumption:

1. Conditional independence becomes finite mixture: Henry et al. (2014) and more.
2. The infeasible moment function  $m(Y_i(1), Y_i(0); \theta)$  becomes linear in feasible moment functions; a plug-in GMM estimator for DTE parameters.

For continuous  $U_i$ , sieve MLE and semiparametric estimation theory:

Shen (1997); Chen and Shen (1998); Ai and Chen (2003) and more. sieve

Need strong assumptions; e.g. bounded support of  $Y_i$  and  $X_i$ .

Why? DTE parameters are complex nonlinear functionals of conditional densities.

## Implementation: finite mixture representation

Under **Assumptions 1-2** that  $Y_i(1), Y_i(0), X_i, (D_i, Z_i) \mid U_i \sim \text{ind}$ ,  
conditional density of  $(Y_i, X_i)$  given  $(D_i, Z_i)$  admits a mixture interpretation:

$$\mathbf{H}_d(y) = \Gamma_X \cdot \Delta_d(y) \cdot \Lambda_d$$
$$f_{Y,X|D,Z}(y, x|d, z) = \int_{\mathbb{R}} \underbrace{f_{X|U}(x|u) \cdot f_{Y(d)|U}(y|u)}_{=\text{mixture component density}} \cdot \underbrace{f_{U|D,Z}(u|d, z)}_{=\text{mixture weights}} du.$$

A change in  $z$  only shifts mixture weights  $f_{U|D,Z}(\cdot|d, z)$ ,  
keeping the same mixture component densities  $\{f_{Y(d),X}(\cdot, \cdot|u)\}_u$ .

A finite support assumption on  $U_i$  gives us a finite mixture.

$K$  is the number of mixture components.

## Implementation: finite mixture representation

With some  $M \geq K$ , construct a  $K$ -way partition for  $Z_i$  and a  $M$ -way partition for  $(Y_i, X_i)$ .

$$\mathbf{H}_d = \begin{pmatrix} \Pr \{Y_i = y^1, X_i = x^1 | D_i = d, Z_i = z^1\} & \cdots & \Pr \{Y_i = y^1, X_i = x^1 | D_i = d, Z_i = z^K\} \\ \vdots & \ddots & \vdots \\ \Pr \{Y_i = y^M, X_i = x^M | D_i = d, Z_i = z^1\} & \cdots & \Pr \{Y_i = y^M, X_i = x^M | D_i = d, Z_i = z^K\} \end{pmatrix}.$$

Then,

$$\underbrace{\mathbf{H}_d}_{M \times K} = \underbrace{\Gamma_d}_{M \times K} \cdot \underbrace{\Lambda_d}_{K \times K}$$

$$\text{with } \Gamma_d = \left( \Pr \{Y_i(d) = y^m, X_i = x^m | U_i = u^k\} \right)_{m,k} \quad (\text{mixture component distributions})$$

$$\Lambda_d = \left( \Pr \{U_i = u^k | D_i = d, Z_i = z^j\} \right)_{k,j}. \quad (\text{mixture weights across subpopulations})$$

Theorem 1 says  $\Gamma_d$  and  $\Lambda_d$  are identified, up to some relabeling.

## Implementation

A GMM model for DTE parameter  $\theta$  with

$$\mathbf{E} [m(Y_i(1), Y_i(0); \theta)] = 0.$$

**Step 0.** Construct a feasible moment function  $\tilde{m}$  with  $\Lambda_d$  as nuisance parameter.

**Step 1.** Estimate  $\Lambda_d = \left( f_{U|D,Z}(u|d, z) \right)_{u,z}$ .

- Estimate the finite mixture model using nonnegative matrix factorization  
 $\Leftrightarrow$  decompose  $\mathbf{H}_d = \Gamma_d \cdot \Lambda_d$  for  $d = 0, 1$ .

**Step 2.** Plug-in GMM to estimate  $\theta$ .

## Implementation, step 0: feasible GMM

Conditional independence/finite mixture  $\mathbf{H}_d = \Gamma_d \cdot \Lambda_d$  implies

$$\begin{aligned} & \left( f_{Y|D,Z}(y|d, z^1) \quad \cdots \quad f_{Y|D,Z}(y|d, z^K) \right) \\ &= \left( f_{Y(d)|U}(y|u^1) \quad \cdots \quad f_{Y(d)|U}(y|u^K) \right) \\ & \quad \cdot \underbrace{\begin{pmatrix} \Pr\{U_i = u^1|D_i = d, Z_i = z^1\} & \cdots & \Pr\{U_i = u^1|D_i = d, Z_i = z^K\} \\ \vdots & \ddots & \vdots \\ \Pr\{U_i = u^K|D_i = d, Z_i = z^1\} & \cdots & \Pr\{U_i = u^K|D_i = d, Z_i = z^K\} \end{pmatrix}}_{=\Lambda_d}. \end{aligned}$$

Let  $\tilde{\lambda}_{jk,d}$  be the  $j$ -th row,  $k$ -th column entry of  $\tilde{\Lambda}_d = (\Lambda_d)^{-1}$ . Then,

$$f_{Y(d)|U}(y|u^k) = \sum_{j=1}^K \tilde{\lambda}_{jk,d} f_{Y|D,Z}(y|d, z^j).$$



## Implementation, step 0: feasible GMM

By substituting for  $f_{Y(d)|U}$ ,

$$\begin{aligned} 0 &= \mathbf{E} [m(Y_i(1), Y_i(0); \theta)] \\ &= \int_{\mathbb{R}} \int_{\mathbb{R}} m(y', y; \theta) f_{Y(0), Y(1)}(y, y') dy dy' \\ &= \int_{\mathbb{R}} \int_{\mathbb{R}} m(y', y; \theta) \left( \sum_{k=1}^K p_U(k) \cdot f_{Y(0)|U}(y|u) \cdot f_{Y(1)|U}(y'|u) \right) dy dy' \\ &= \int_{\mathbb{R}} \int_{\mathbb{R}} m(y', y; \theta) \left( \sum_{k=1}^K p_U(k) \cdot \sum_{j=1}^K \tilde{\lambda}_{jk,0} f_{Y|D=0,Z}(y|z^j) \cdot \sum_{j'=1}^K \tilde{\lambda}_{j'k,1} f_{Y|D=1,Z}(y|z^{j'}) \right) dy dy' \\ &= \sum_{k=1}^K \sum_{j=1}^K \sum_{j'=1}^K p_U(k) \tilde{\lambda}_{jk,0} \tilde{\lambda}_{j'k,1} \underbrace{\int_{\mathbb{R}} \int_{\mathbb{R}} m(y', y; \theta) \left( f_{Y|D=0,Z}(y|z^j) \cdot f_{Y|D=1,Z}(y|z^{j'}) \right) dy dy'}_{\text{quadratic moment of } (Y_i, D_i, X_i, Z_i)} \end{aligned}$$

where  $p_U(k) = \Pr \{U_i = u^k\}$

## Implementation, step 0: feasible GMM

For example,

$$\begin{aligned} F_{Y(1)-Y(0)}(\delta) &= \mathbf{E}[\mathbf{1}\{Y_i(1) \leq Y_i(0) + \delta\}] \\ &= \sum_{k=1}^K \sum_{j=1}^K \sum_{j'=1}^K \frac{p_U(k) \tilde{\lambda}_{jk,1} \tilde{\lambda}_{j'k,0}}{p_{D,Z}(1,j) p_{D,Z}(0,j')} \\ &\quad \cdot \mathbf{E}[\mathbf{1}\{Y_i \leq Y_{i'} + \delta, D_i = 1, Z_i = z^j, D_{i'} = 0, Z_{i'} = z^{j'}\}] \end{aligned}$$

for all  $\delta \in \mathbb{R}$ , with  $(Y_i, D_i, Z_i) \perp\!\!\!\perp (Y_{i'}, D_{i'}, Z_{i'})$ . derivation

The nuisance parameters are  $\tilde{\lambda} = \text{vec}\left((\Lambda_0)^{-1}, (\Lambda_1)^{-1}\right)$  and

$$p = \text{vec}\left(\{p_{D,U}(d,k)\}_{d,k}, \{p_{D,Z}(d,j)\}_{d,j}\right).$$

$\tilde{\lambda}$  is matrix inverses of  $(\Lambda_0, \Lambda_1)$ .

$p$  collects joint probabilities of  $(D_i, U_i)$  and  $(D_i, Z_i)$ ; also a function of  $(\Lambda_0, \Lambda_1)$ .

**Proposition 1.**

Assumptions 1-3 hold. Suppose that a parameter of interest  $\theta$  is identified by an infeasible moment condition

$$\mathbf{E} \left[ m(Y_i(1), Y_i(0), D_i, X_i; \theta) \right] = 0.$$

Then, there is a moment function  $\tilde{m}$  such that  $\theta$  is identified by a feasible quadratic moment condition

$$\mathbf{E} \left[ \tilde{m} \left( (Y_i, D_i, X_i, Z_i), (Y_{i'}, D_{i'}, X_{i'}, Z_{i'}); \theta, \tilde{\lambda}, p \right) \right] = 0$$

where  $i \neq i'$ .

The proxy  $Z_i$  is used to shift  $U_i$ ;

for DTE parameters that involve  $Z_i$ , a similar result with higher-order moments.

## Implementation, step 1: nonnegative matrix factorization

For first-step nuisance parameter estimation,  
use **nonnegative matrix factorization** to estimate  $(\Lambda_0, \Lambda_1)$ .

Recall  $\mathbf{H}_d = \Gamma_d \cdot \Lambda_d$ .

Given discretized  $Y_i, X_i$  and  $Z_i$ , estimate  $\mathbf{H}_d$  with sample analogue:

$$\mathbb{H}_d = \left( \frac{\sum_{i=1}^n \mathbf{1}\{Y_i = y^m, D_i = d, X_i = x^m, Z_i = z^k\}}{\sum_{i=1}^n \mathbf{1}\{D_i = d, Z_i = z^k\}} \right)_{m,k}$$

-  $\mathbb{H}_d$  is a  $\sqrt{n}$ -consistent estimator of  $\mathbf{H}_d$ .

## Implementation, step 1: nonnegative matrix factorization

Solve the following nonnegative matrix factorization problem:

$$\left(\hat{\Gamma}_0, \hat{\Gamma}_1, \hat{\Lambda}_0, \hat{\Lambda}_1\right) = \arg \min \left\|\mathbb{H}_0 - \Gamma_0 \cdot \Lambda_0\right\|_F + \left\|\mathbb{H}_1 - \Gamma_1 \cdot \Lambda_1\right\|_F \quad (1)$$

subject to 1)  $\Gamma_0, \Gamma_1, \Lambda_0, \Lambda_1$  are nonnegative.

Also, their columnwise sums are one.  $\dots$  (*linear constraints*)

2)  $\Gamma_0$  and  $\Gamma_1$  satisfy  $Y_i(d) \perp\!\!\!\perp X_i \mid U_i \dots$  (*quadratic constraints*)

3)  $\Gamma_0$  and  $\Gamma_1$  satisfy  $D_i \perp\!\!\!\perp X_i \mid U_i \dots$  (*linear constraints*)

This optimization is principal component analysis + additional constraint. v. PCA

(1) is solved iteratively; higher computational burden. algorithm

## Implementation, step 1: nonnegative matrix factorization

The first-step nuisance parameter estimators are:

$$\hat{\lambda} = \text{vec}\left(\left(\hat{\Lambda}_0\right)^{-1}, \left(\hat{\Lambda}_1\right)^{-1}\right) \quad \text{and} \quad \hat{p} = \text{vec}\left(\{\hat{p}_{D,U}(d, k)\}_{d,k}, \{\hat{p}_{D,Z}(d, j)\}_{d,j}\right)$$

where

$$\hat{p}_{D,Z}(d, j) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{D_i = d, Z_i = z^j\} \quad \forall d = 0, 1 \text{ and } j = 1, \dots, K$$

$$\begin{pmatrix} \hat{p}_{D,U}(d, 1) \\ \vdots \\ \hat{p}_{D,U}(d, K) \end{pmatrix} = \hat{\Lambda}_d \begin{pmatrix} \hat{p}_{D,Z}(d, 1) \\ \vdots \\ \hat{p}_{D,Z}(d, K) \end{pmatrix} \quad \forall d = 0, 1$$

since  $p_{D,U}(d, k) = \sum_{j=1}^K \Pr\{D_i = d, U_i = u^k, Z_i = z^j\}$

$$= \underbrace{\sum_{j=1}^K \Pr\{U_i = u^k | D_i = d, Z_i = z^j\}}_{=\text{components of the } k\text{-th row of } \Lambda_d} \cdot \Pr\{D_i = d, Z_i = z^j\}.$$

## Implementation, step 1: nonnegative matrix factorization

**Theorem 2.**

Assumptions 1-3 hold. Up to some permutation on  $\{u^1, \dots, u^K\}$ ,

$$\left\| \hat{\Lambda}_0 - \Lambda_0 \right\|_F = O_p \left( \frac{1}{\sqrt{n}} \right) \quad \text{and} \quad \left\| \hat{\Lambda}_1 - \Lambda_1 \right\|_F = O_p \left( \frac{1}{\sqrt{n}} \right)$$

as  $n \rightarrow \infty$ .

A direct corollary is that  $\hat{\lambda}$  and  $\hat{p}$  are consistent for  $\tilde{\lambda}$  and  $p$  at the same rate.

The same rate can be established for  $\hat{\Gamma}_0$  and  $\hat{\Gamma}_1$ .

## Implementation: comparison to diagonalization

Existing estimators for nonparametric finite mixture rely on (joint) diagonalization.

The same  $n^{-\frac{1}{2}}$  rate and asymptotic normality.

Hu (2008); Kasahara and Shimotsu (2009); Bonhomme et al. (2016) and more.

$$\mathbb{H}_d(y) \left( \sum_{y'} \mathbb{H}_d(y') \right)^{-1} = \hat{\Gamma}_X \cdot \hat{\Delta}_d(y) \cdot \left( \hat{\Gamma}_X \right)^{-1}.$$

Additional regularization in the NMF estimator:

- Nonnegativity constraints for  $\hat{f}_{X|U}$ . (*eigenvectors*)
- Sum-to-one constraints for  $\hat{f}_{Y(d)|U}$ . (*eigenvalues*)

DTE estimation depends on nuisance parameter estimation:  $(\Lambda_0)^{-1}$  and  $(\Lambda_1)^{-1}$ .

Additional regularization helps.



## Implementation, step 2: Neyman orthogonality

Nuisance parameter estimation has first-order impact on DTE estimation.

Thus, I orthogonalize the feasible moment function  $\tilde{m}$ . Delta method

$\hat{\lambda}$  and  $\hat{p}$  are highly nonlinear functions of  $\mathbb{H}_0$  and  $\mathbb{H}_1$ .

No usual “first-order condition”-type moments available.

1. For  $\tilde{\lambda}$ , use the quadratic moments of conditional independence: for  $d = 0, 1$ ,

$$\Pr \{Y_i(d) = y, X_i = x | U_i = u\} = \Pr \{Y_i(d) = y | U_i = u\} \cdot \Pr \{X_i = x | U_i = u\}.$$

more

2. For  $p_{D,U}$ , use the moments of law of iterated expectation:

$$\Pr \{D_i = d, X_i = x\} = \sum_{k=1}^k p_{D,U}(d, k) \Pr \{X_i = x | U_i = u^k\}.$$

## Implementation, step 2: Neyman orthogonality

Let  $\phi$  be the moment function for the additional moments.

**Lemma 1.** Assumptions 1-3 hold. Then,  $\begin{pmatrix} \mathbf{E}[\frac{\partial}{\partial \lambda} \phi] \\ \mathbf{E}[\frac{\partial}{\partial p} \phi] \end{pmatrix}$  has full row rank.

With  $\mu = \begin{pmatrix} \mathbf{E}[\frac{\partial}{\partial \lambda} \phi] \\ \mathbf{E}[\frac{\partial}{\partial p} \phi] \end{pmatrix}^+ \begin{pmatrix} \mathbf{E}[\frac{\partial}{\partial \lambda} \tilde{m}] \\ \mathbf{E}[\frac{\partial}{\partial p} \tilde{m}] \end{pmatrix}$ , the **orthogonalized moment function** is

$$\begin{aligned} & \psi \left( (Y_i, D_i, X_i, Z_i), (Y_{i'}, D_{i'}, X_{i'}, Z_{i'}); \theta, \tilde{\lambda}, p, \mu \right) \\ &= \tilde{m} \left( (Y_i, D_i, X_i, Z_i), (Y_{i'}, D_{i'}, X_{i'}, Z_{i'}); \theta, \tilde{\lambda}, p \right) - \mu^\top \phi \left( (Y_i, D_i, X_i, Z_i), (Y_{i'}, D_{i'}, X_{i'}, Z_{i'}); \tilde{\lambda}, p \right) \end{aligned}$$

This applies to any moment-identified parameter in a finite mixture model.

## Implementation, step 2: plug-in GMM

Using Neyman orthogonality, asymptotic normality is established.

### **Theorem 3.**

Assumptions 1-3 hold.  $\hat{\theta}$  is the plug-in GMM estimator of  $\theta$ , using  $\psi$ . Then,

$$\sqrt{n} \left( \hat{\theta} - \theta \right) \xrightarrow{d} \mathcal{N}(0, \Sigma)$$

as  $n \rightarrow \infty$  with some consistently estimable  $\Sigma$ .

Inference is computationally easy;

asymptotic standard error is computed from  $\psi \left( (Y_i, D_i, X_i, Z_i), (Y_{i'}, D_{i'}, X_{i'}, Z_{i'}); \hat{\theta}, \hat{\lambda}, \hat{p}, \hat{\mu} \right)$ .

Wide applicability beyond DTE estimation:

random coefficient model, dynamic discrete choice and more.

## Implementation: falsification test

Firstly, we can test  $X_i \perp\!\!\!\perp D_i \mid U_i$  from Assumption 1 as a falsification test:

$$\sum_{k=1}^K \sum_{m=1}^{M_X} \left( f_{X|D=1,U}(x^m|u^k) - f_{X|D=0,U}(x^m|u^k) \right)^2 = 0.$$

Alternatively, we can test  $D_i \perp\!\!\!\perp U_i$  when  $D_i$  is randomly assigned:

$$\sum_{k=1}^K \left( p_{U|D=1}(u^k) - p_{U|D=0}(u^k) \right)^2 = 0.$$

### Theorem 4.

Under Assumptions 1-3,  $T_n^1 \xrightarrow{d} \chi^2(K \cdot M_X)$  as  $n \rightarrow \infty$ .

Additionally, when  $D_i \perp\!\!\!\perp U_i$ ,  $T_n^2 \xrightarrow{d} \chi^2(K)$  as  $n \rightarrow \infty$ . degrees of freedom

Simulation

## Simulation

Monte Carlo simulations ( $B = 1000$ ) with two DGPs where  $X_i, Z_i, U_i \in \{1, 2, 3\}$ .

$$Y_i(d) \mid (U_i = k) \sim \mathcal{N}(\mu^k(d), \sigma^k(d)^2).$$

and

$$\left( \Pr\{X_i = x \mid U_i = k\} \right)_{x,k} = \begin{pmatrix} 0.911 & 0.050 & 0.022 \\ 0.067 & 0.900 & 0.067 \\ 0.022 & 0.050 & 0.911 \end{pmatrix}$$
$$\left( \Pr\{Z_i = z \mid U_i = k\} \right)_{z,k} = \begin{pmatrix} 0.689 & 0.175 & 0.078 \\ 0.233 & 0.650 & 0.233 \\ 0.078 & 0.175 & 0.689 \end{pmatrix}, \quad \begin{pmatrix} 0.911 & 0.050 & 0.022 \\ 0.067 & 0.900 & 0.067 \\ 0.022 & 0.050 & 0.911 \end{pmatrix}.$$

The smallest singular values for  $\Lambda$  is 0.337 and 0.806. [specifics](#)

Estimate marginal distribution of treatment effect  $F_{Y(1)-Y(0)}(\delta)$ .

	true value	bias				rMSE			
$\hat{F}_{Y(1)-Y(0)}(0)$	0.084	0.000	0.000	-0.002	-0.001	0.014	0.009	0.011	0.007
$\hat{F}_{Y(1)-Y(0)}(1)$	0.264	0.001	0.001	-0.001	-0.001	0.023	0.015	0.019	0.012
$\hat{F}_{Y(1)-Y(0)}(2)$	0.536	0.001	0.000	0.000	-0.001	0.025	0.016	0.022	0.014
$\hat{F}_{Y(1)-Y(0)}(3)$	0.775	0.002	0.000	0.002	0.000	0.020	0.012	0.018	0.011
$\hat{F}_{Y(1)-Y(0)}(4)$	0.911	0.005	0.002	0.003	0.001	0.014	0.008	0.012	0.007
$\sigma_{\min}(\Lambda)$		0.337	0.337	0.806	0.806	0.337	0.337	0.806	0.806
$n$		750	2000	750	2000	750	2000	750	2000

Table 1: Bias and rMSE of DTE estimator  $\hat{F}_{Y(1)-Y(0)}(\delta)$  based on NMF.

Estimation performance improves as  $Z_i$  gets more informative, i.e.  $\sigma_{\min}(\Lambda)$  goes up.

	true value	coverage probability			
$\hat{F}_{Y(1)-Y(0)}(0)$	0.084	0.971	0.951	0.952	0.935
$\hat{F}_{Y(1)-Y(0)}(1)$	0.264	0.975	0.959	0.958	0.952
$\hat{F}_{Y(1)-Y(0)}(2)$	0.536	0.970	0.960	0.957	0.951
$\hat{F}_{Y(1)-Y(0)}(3)$	0.775	0.962	0.959	0.943	0.951
$\hat{F}_{Y(1)-Y(0)}(4)$	0.911	0.940	0.954	0.934	0.948
$\sigma_{\min}(\Lambda)$		0.337	0.337	0.806	0.806
$n$		750	2000	750	2000

Table 2: Coverage of 95% confidence interval based on NMF.

Slight conservatism when  $\sigma_{\min}(\Lambda)$  is low and  $n$  is small.



	true value	bias				rMSE			
$\hat{F}_{Y(1)-Y(0)}(0)$	0.084	0.000	0.000	0.014	0.008	0.014	0.009	0.034	0.029
$\hat{F}_{Y(1)-Y(0)}(1)$	0.264	0.001	0.001	0.006	0.004	0.023	0.015	0.030	0.021
$\hat{F}_{Y(1)-Y(0)}(2)$	0.536	0.001	0.000	-0.006	-0.005	0.025	0.016	0.037	0.029
$\hat{F}_{Y(1)-Y(0)}(3)$	0.775	0.002	0.000	-0.009	-0.007	0.020	0.012	0.040	0.032
$\hat{F}_{Y(1)-Y(0)}(4)$	0.911	0.005	0.002	-0.006	-0.004	0.014	0.008	0.025	0.019
first-step		NMF	NMF	EVD	EVD	NMF	NMF	EVD	EVD
$n$		750	2000	750	2000	750	2000	750	2000

Table 3: Comparison between first-step NMF and EVD, when  $\sigma_{\min}(\Lambda) = 0.337$ .

For EVD, we get nonzero bias and rMSE 1.25-4.77 times larger: intensive margin.

	success rate				computation time (sec)			
NMF	0.999	1.000	1.000	1.000	98.01	163.28	66.32	117.40
EVD	0.528	0.666	0.790	0.846	19.27	80.57	19.57	73.77
$\sigma_{\min}(\Lambda)$	0.337	0.337	0.806	0.806	0.337	0.337	0.806	0.806
$n$	750	2000	750	2000	750	2000	750	2000

Table 4: Success rate and computation time for DTE estimation based on NMF and EVD.

The estimation halted for 15.4-47.2% of the samples: extensive margin.

## Empirical Illustration

## Empirical illustration: setup

I revisit Jones et al. (2019), which studies the effect of workplace wellness program. The program *eligibility* was randomly assigned to employees at UIUC; intent-to-treat.

The variables in the dataset are:

$Y_i$  = monthly medical spending over August 2016-July 2017

$D_i = \mathbf{1}\{\text{eligible for the wellness program starting in September 2016}\}$

$X_i$  = monthly medical spending over July 2015-July 2016

$Z_i$  = monthly medical spending over August 2017-January 2019

*“This year’s underlying health status  $U_i$  only depends on last year’s health status.”*

## Empirical illustration: choice of $K$

1. Consider a  $M_X \times 2M_Z$  matrix  $\mathbf{H}_X$ :

$$\mathbf{H}_X = \begin{pmatrix} \Pr\{X_i = x^1 | (D_i, Z_i) = (0, z^1)\} & \cdots & \Pr\{X_i = x^1 | (D_i, Z_i) = (1, z^{M_Z})\} \\ \vdots & \ddots & \vdots \\ \Pr\{X_i = x^{M_X} | (D_i, Z_i) = (0, z^1)\} & \cdots & \Pr\{X_i = x^{M_X} | (D_i, Z_i) = (1, z^{M_Z})\} \end{pmatrix}.$$

$\mathbf{H}_X$  pools information from  $\{i : D_i = 0\}$  and  $\{i : D_i = 1\}$  and should be at most rank  $K$ .

Apply eigenvalue ratio estimator and rank test.

2. With true  $K$ , estimated densities should satisfy

$$f_{X|D=1,U}(x|u) = f_{X|D=0,U}(x,u) \quad \forall x, u,$$

$$f_{U|D=1}(u) = f_{U|D=0}(u) \quad \forall u.$$

Apply falsification tests.

## Empirical illustration: choice of $K$

Both rank test and eigenvalue ratio estimator suggest  $K = 3$ .

$K$	1	2	3	4	5	6	7	8
eigenvalue ratio	3.505	3.991	4.029	2.721	1.653	1.863	1.418	3.309
growth ratio	0.964	1.135	1.472	1.353	0.893	0.956	0.580	1.035

Table 5: Eigenvalue ratios and growth ratios

$K$	1	2	3	4	5	6
test statistic	884.82	116.23	35.75	20.08	13.80	7.94
$p$ -value	0.000	0.001	0.984	0.998	0.995	0.992

Table 6: Kleibergen-Paap rank test statistics for  $H_0 : \text{rank} = K$  and their  $p$ -values

## Empirical illustration: choice of $K$

Two falsification test statistics:

$T_n^1 = \chi^2$  test statistic for  $f_{X|D=1,U}(x|u) = f_{X|D=0,U}(x,u) \quad \forall x, u,$

$T_n^2 = \chi^2$  test statistic for  $f_{U|D=1}(u) = f_{U|D=0}(u) \quad \forall u.$

$K$	3	4	5	6
$T_n^1$	17.68	27.07	16.79	47.66
$p$ -value	0.477	0.301	0.975	0.092
$T_n^2$	1.57	0.22	0.24	4.27
$p$ -value	0.666	0.995	0.999	0.640

Table 7: Falsification test statistics  $(T_n^1, T_n^2)$  and their  $p$ -values

## Empirical illustration: conditional ATE

	(1)	(2)	(3)	(4)	(5)
CATE( $\mathcal{Y}$ ) (\$)	-20.43 (166.31)	164.17 (389.05)	297.49 (350.17)	-157.24 (515.51)	-732.94* (443.00)
$\mathcal{Y}$	[0, 42.7]	(42.7, 132.2]	(132.2, 286.8]	(286.8, 671.1]	(671.1, $\infty$ ]

Table 8: Conditional average treatment effect  $\mathbb{E}[Y_i(1) - Y_i(0)|Y_i(0) \in \mathcal{Y}]$ .

Conditional ATE across five quintiles of  $Y_i(0)$ :

$$\mathcal{Y} = (F_{Y(0)}^{-1}(0), F_{Y(0)}^{-1}(1/5)], \dots, (F_{Y(0)}^{-1}(4/5), F_{Y(0)}^{-1}(1)].$$

In Jones et al. (2019),  $p$ -values for ATE are 0.86-0.94. On page 1890 of Jones et al. (2019),

*“There may exist subpopulations who did benefit from the intervention or who would have benefited had they participated.”*



## Empirical illustration: treatment effect distribution

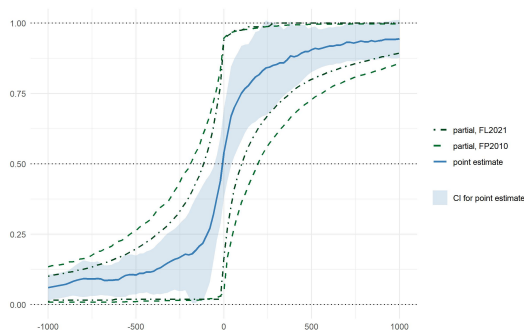


Figure 2: Marginal distribution of  $Y_i(1) - Y_i(0)$ .

Fan and Park (2010):  $(Y_i(1), Y_i(0)) \perp\!\!\!\perp D_i$ .

Frandsen and Lefgren (2021):  $(Y_i(1), Y_i(0)) \perp\!\!\!\perp D_i$ ,

$F_{Y(1)|Y(0)}(y|y')$  is decreasing in  $y'$  for all  $y$ ,

$F_{Y(0)|Y(1)}(y|y')$  is decreasing in  $y'$  for all  $y$ .

## Empirical illustration: treatment effect distribution

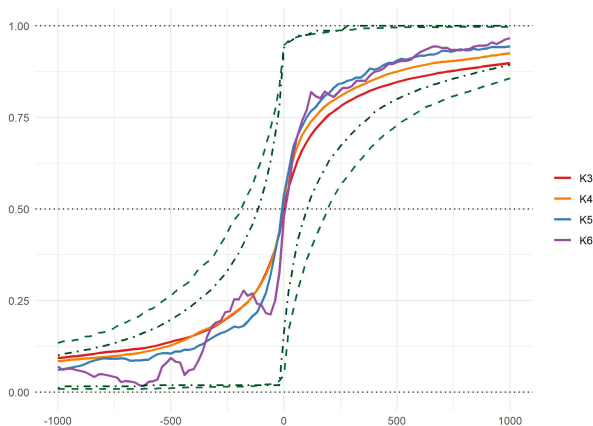
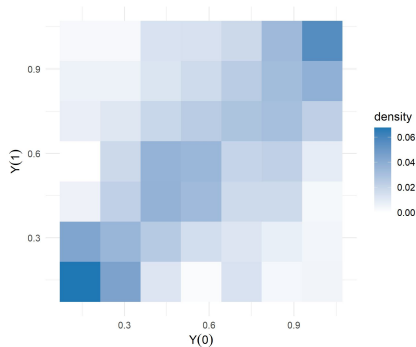


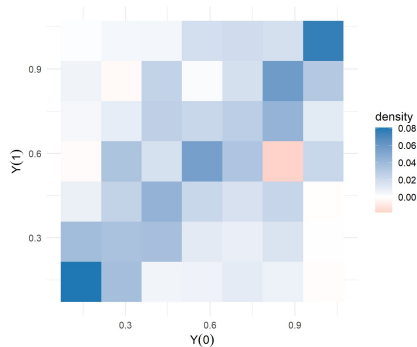
Figure 3: Marginal distribution of  $Y_i(1) - Y_i(0)$ , across  $K$ .

For 37% of the support,  $\hat{F}_{Y(1)-Y(0)}$  with  $K = 6$  was decreasing.  
Possible misspecification/discretization bias for on the right tail.

## Empirical illustration: joint density of $Y_i(1)$ and $Y_i(0)$



(a)  $K = 4$



(b)  $K = 5$

Figure 4: Joint density of  $Y_i(1)$  and  $Y_i(0)$ , across  $K = 4, 5$ .

High correlation on the two ends of the spectrum.

## Conclusion

- Assume a latent variable  $U$  such that

$$Y_i(1) \perp\!\!\!\perp Y_i(0) \mid U_i.$$

This assumption could be thought of as a ‘latent rank invariance’ condition.

- Use two proxy variables  $X_i$  and  $Z_i$  to identify the distribution of  $Y_i(d)|U_i$ .  
Measurement error model and hidden Markov model motivate the use of proxies.
- Nonnegative matrix factorization estimates finite mixture.  
Better finite sample performance due to additional regularization.
- An asymptotic distribution is derived for the DTE estimator.

## Conditional independence framework

**Assumption 1.**  $(Y_i(1), Y_i(0), X_i) \perp\!\!\!\perp (D_i, Z_i) \mid U_i$ .

- In proximal inference literature,

$X_i$  = outcome-aligned proxy and  $Z_i$  = treatment-aligned proxy.

Miao et al. (2018); Deaner (2023); Nagasawa (2022) and more.

- Distribution of  $(Y_i(d), X_i) \mid U_i$  is set identified: Henry et al. (2014).

**Assumption 2.**  $Y_i(1), Y_i(0), X_i$  are mutually independent given  $U_i$ .

- Nonclassical measurement error literature:  $Y_i, X_i, Z_i \mid U_i \sim \text{ind.}$

Hu (2008); Hu and Schennach (2008) and more.

- Conditional distributions are point identified.

- Extended to a potential outcome setup.

Additionally,  $Y_i(1) \perp\!\!\!\perp Y_i(0) \mid U_i$  recovers the joint dist. of  $Y_i(1)$  and  $Y_i(0)$ .

## Conditional independence: example 1 (*repeated measurement*)

Attanasio et al. (2020): early childhood intervention's effect on children's development.

$Y_i$  is test score at follow-up,  $U_i$  is innate ability at baseline, and  $(X_i, Z_i)$  are test scores at baseline.

$$Y_i(d) = \mu^d + \alpha^d U_i + \varepsilon_i^d \quad \text{for } d = 0, 1,$$

$$X_i = \mu^X + \alpha^X U_i + \varepsilon_i^X,$$

$$Z_i = \mu^Z + \alpha^Z U_i + \varepsilon_i^Z.$$

## Conditional independence: example 1 (*repeated measurement*)

Attanasio et al. (2020): early childhood intervention's effect on children's development.

$Y_i$  is test score at follow-up,  $U_i$  is innate ability at baseline, and  $(X_i, Z_i)$  are test scores at baseline.

$$Y_i(d) = \mu^d + \alpha^d U_i + \varepsilon_i^d \quad \text{for } d = 0, 1,$$

$$X_i = \mu^X + \alpha^X U_i + \varepsilon_i^X,$$

$$Z_i = \mu^Z + \alpha^Z U_i + \varepsilon_i^Z.$$

Assumptions 1-2 hold when

- $\varepsilon_i^0, \varepsilon_i^1, \varepsilon_i^X$  and  $\varepsilon_i^Z$  are mutually independent given  $U_i$ .
- $D_i$  is randomly assigned.

## Conditional independence: example 2 (*past and future outcomes*)

A common shock process  $\{V_{it}\}_{t=1}^3$  and random shocks  $(\varepsilon_{i1}^0, \varepsilon_{i2}^0, \varepsilon_{i2}^1, \varepsilon_{i3}^0, \varepsilon_{i3}^1)$ .

$$Y_{it}(d) = g_d(V_{it}, \varepsilon_{it}^d) \quad \text{for } d = 0, 1 \text{ and } t = 1, 2, 3,$$
$$Y_{it} = \begin{cases} Y_{i1}(0) & \text{if } t = 1 \\ D_i \cdot Y_{it}(1) + (1 - D_i) \cdot Y_{it}(0) & \text{if } t \geq 2 \end{cases}.$$



## Conditional independence: example 2 (*past and future outcomes*)

A common shock process  $\{V_{it}\}_{t=1}^3$  and random shocks  $(\varepsilon_{i1}^0, \varepsilon_{i2}^0, \varepsilon_{i2}^1, \varepsilon_{i3}^0, \varepsilon_{i3}^1)$ .

$$Y_{it}(d) = g_d(V_{it}, \varepsilon_{it}^d) \quad \text{for } d = 0, 1 \text{ and } t = 1, 2, 3,$$
$$Y_{it} = \begin{cases} Y_{i1}(0) & \text{if } t = 1 \\ D_i \cdot Y_{it}(1) + (1 - D_i) \cdot Y_{it}(0) & \text{if } t \geq 2 \end{cases}.$$

Assumptions 1-2 hold with  $(Y_i, X_i, Z_i, U_i) = (Y_{i2}, Y_{i1}, Y_{i3}, V_{i2})$  when

- $(\{V_{it}\}_{t=1}^3, D_i), \varepsilon_{i1}^0, \varepsilon_{i2}^0, \varepsilon_{i2}^1, \varepsilon_{i3}^0, \varepsilon_{i3}^1$  are mutually independent.
- $\{V_{it}\}_{t=1}^3$  is first-order Markovian given  $D_i$ :  $V_{i1} \perp\!\!\!\perp V_{i3} \mid (V_{i2}, D_i)$ .
- $D_i$  is randomly assigned at time  $t = 2$ :  $\{V_{it}\}_{t=1}^2 \perp\!\!\!\perp D_i$ .

## Regime-changing treatment effect mechanism

In both measurement error model and hidden Markov model,

- Two separate outcome generating processes for  $Y_i(1)$  and  $Y_i(0)$ : *regime-changing*.

$$Y_i(0) = g_0(U_i, \varepsilon_i^0),$$

$$Y_i(1) = g_1(U_i, \varepsilon_i^1).$$

The regime-specific random shocks are purely random, satisfying  $\varepsilon_i^1 \perp\!\!\!\perp \varepsilon_i^0 \mid U_i$ .

In contrast to *input-changing* treatment mechanism. [more](#)

## Examples of regime-changing treatment mechanism

Thus, Assumption 2 is most plausible when the treatment induces systemic changes:

- Attanasio et al. (2020):  
Treatment provided parenting guidance, changing how parents interacted with children.
- Jones et al. (2019): ← my empirical example  
Treatment provided information sessions on healthy lifestyle, changing how participants sought medical service and took self-care measures.
- Job assignment: e.g. the National Supported Work Demonstration.  
Treatment changes how worker skill  $U_i$  leads to outcome  $Y_i$  such as income.
- Teaching methodology: e.g. Banerjee et al. (2007); Muralidharan et al. (2019).  
Treatment changes how student aptitude  $U_i$  leads to outcome  $Y_i$  such as academic achievement.

## Input-changing treatment mechanism

Two common independence assumptions:

$$Y_i(1) \perp\!\!\!\perp Y_i(0) \quad \text{and} \quad Y_i(0) \perp\!\!\!\perp (Y_i(1) - Y_i(0))$$

The latter can be motivated by *input-changing* treatment mechanism: with  $V_i \perp\!\!\!\perp \varepsilon_i \mid U_i$ ,

$$Y_i(d) = \alpha + \mu^0 U_i + d \cdot \mu^1 V_i + \varepsilon_i$$

satisfies  $Y_i(0) \perp\!\!\!\perp (Y_i(1) - Y_i(0)) \mid U_i$ .

Treatment turns on a new source of individual-level heterogeneity  $V_i$ , which is (conditionally) independent of the existing heterogeneity  $\varepsilon_i$ .

For example, providing new infrastructure: computers in teaching environment. [back](#)

## Spectral Theorem of Hu and Schennach (2008)

A linear operator  $L_{Y=y, X|D=d, X}$  maps a density of  $Z_i$  to a density of  $(Y_i(d) = y, X_i)$ :

$$(L_{Y=y, X|D=d, Z} g)(x) = \int_{\mathbb{R}} f_{Y(d), X|D, Z}(y, x|d, z) g(z) dz.$$

From the decomposition based on Assumption 2, we get

$$L_{Y=y, X|D=d, Z} = L_{X|U} \cdot \Delta_{Y=y|U} \cdot L_{U|D=d, Z}$$

with similarly defined operators  $L_{X|U}$ ,  $L_{U|D=d, Z}$  and a diagonal operator  $\Delta_{Y=y|U}$ . Thus,

$$\begin{aligned} L_{Y=y, X|D=d, Z} (L_{X|D=d, Z})^{-1} &= L_{X|U} \cdot \Delta_{Y=y|U} \cdot L_{U|D=d, Z} \cdot (L_{X|U} \cdot L_{U|D=d, Z})^{-1} \\ &= \underbrace{L_{X|U} \cdot \Delta_{Y=y|U} \cdot (L_{X|U})^{-1}}_{\text{spectral decomposition}}. \end{aligned}$$

## Assumptions 1-3

### Assumptions 1-2.

$Y_i(1), Y_i(0), X_i, (D_i, Z_i)$  are mutually independent of each other given  $U_i$ .

### Assumption 3.

- a. *(finitely discrete  $U_i$ )*  $U_i \in \{u^1, \dots, u^K\}$ .
- b. *(full rank)*  $\Gamma_X, \Lambda_0$  and  $\Lambda_1$  have rank  $K$ .
- c. *(no repeated eigenvalue)* For any  $k \neq k'$ , there exist some  $d \in \{0, 1\}$  and  $y$  such that

$$\Pr \left\{ Y_i(d) = y | U_i = u^k \right\} \neq \Pr \left\{ Y_i(d) = y | U_i = u^{k'} \right\}.$$

A3

P1

T2

## Assumption 4

### Assumption 4.

- a. (continuous  $U_i$ )  $U_i \in [0, 1]$ .
- b. (bounded density) All marginal and conditional densities of  $(Y_i(1), Y_i(0), X_i, Z_i, U_i)$  are bounded.
- c. (completeness) Let  $f_{X|Z,d}$  denote the conditional density of  $X_i$  given  $(D_i = d, Z_i)$ .

$$\int_{\mathbb{R}} |g(x)| dx \quad \text{and} \quad \int_{\mathbb{R}} g(x) f_{X|Z,d}(x|z) dx = 0 \quad \forall d, z$$

implies  $g(x) = 0$ . Assume similarly for  $f_{X|U}$ .

- d. (no repeated eigenvalue)  $\forall u \neq u'$ , there exists  $d \in \{0, 1\}$  such that

$$\Pr \{ f_{Y(d)|U}(Y_i(d)|u) \neq f_{Y(d)|U}(Y_i(d)|u') | D_i = d \} > 0.$$

## Why do we need Assumption 5?

Under Assumptions 1,2 and 3/4, we have identified

$$\{f_{Y(1)|U}(\cdot|u), f_{Y(0)|U}(\cdot|u), f_{X|U}(\cdot|u), f_{U|D=1,Z}(u|\cdot), f_{U|D=0,Z}(u|\cdot)\}_u.$$

Each group of five is not yet connected to a value of  $u$ ; unordered collection.

When the collection is finite, not having an ordering is okay.

Probability mass function of  $U_i$  is still identified.

When the collection is infinite, not having an ordering matters.

Why? Now we are in density territory.

$\{f_{U|D=1,Z}(u|\cdot), f_{U|D=0,Z}(u|\cdot)\}_u$  gives us  $\{f_U(u)\}_u$  but not  $f_U$ .

Assumption 5 orders  $\{f_U(u)\}_u$  using  $\tilde{u} = \mathbf{E}[Y_i(d)|U_i = u]$ . [back](#)



## Multidimensional $U_i$

Skill formation/human capital accumulation literature often model two-dimensional  $U_i$ :

Carneiro et al. (2003); Cunha et al. (2010); Attanasio et al. (2020) and more

- $U_i = (U_i^C, U_i^N)$ : cognitive and noncognitive ability of a child.
- $X_i = (X_i^C, X_i^N)$ : cognitive ability test scores and noncognitive ability test scores.
- $Z_i = (Z_i^C, Z_i^N)$ : another set of ability scores, measured independently.

Components of  $X_i, Z_i$  need not match components of  $U_i$ .

Helps in assuming that any remaining heterogeneity after controlling for  $U_i$  is purely random.

[back](#)

## Identification: implicit restriction

A crucial step in the identification argument is that there exists some  $w$  such that

$$\mathbf{E}[Y_i(1)|Y_i(0) = y] = \int_{\mathbb{R}} \frac{w(y, z)}{f_{Y(0)}(y)} \cdot \mathbf{E}[Y_i|D_i = 1, Z_i = z]dz,$$
$$\mathbf{E}[Y_i(1)Y_i(0)] = \int_{\mathbb{R}} \int_{\mathbb{R}} w(y, z) \cdot y \mathbf{E}[Y_i|D_i = 1, Z_i = z]dydz.$$

$\mathbf{E}[Y_i|D_i = 1, Z_i]$  replaces  $Y_i(1)$  and  $w(y, z)$  replaces the joint density of  $(Y_i(1), Y_i(0))$ .

“Proxy variable  $Z_i$  creates sufficient variation in the distribution of  $Y_i(1)$ .”

The implicit restriction is that

“conditional distribution of  $Y_i(1)$  given  $Y_i(0)$  is a linear combination of  $\{F_{Y|D=1, Z}(\cdot|z)\}_z$ .”

## Sieve MLE

To allow for a continuous  $U_i$ , we can directly construct a likelihood using sieves:

$$f_{Y,X|D=d,Z,n}(y,x|z;\theta) = \int_{\mathbb{R}} f_{Y(d)|U,n}(y|u;\theta) \cdot f_{X|U,n}(x|u;\theta) \cdot f_{U|D=d,Z,n}(u|z;\theta) du.$$

Nonnegativity, sum-to-one, monotonicity conditions are easy to impose with Bernstein polynomials: a Bernstein polynomial of degree  $m$  is

$$g_m(u) = \sum_{k=0}^m \theta_k u^k (1-u)^{m-k}.$$

Then, monotonicity of  $\int_0^1 u g_m(u) du$  is a set of linear constraints on  $\{\theta_k\}_{k=0}^m$ .

## Choice of $K$

Under Assumption 3, the rank of the following  $M_X \times 2M_Z$  matrix is  $K$ :

$$\mathbf{H}_X = \begin{pmatrix} \Pr \{X_i \in \mathcal{X}^1 | D_i = 0, Z_i \in \mathcal{Z}^1\} & \cdots & \Pr \{X_i \in \mathcal{X}^1 | D_i = 1, Z_i \in \mathcal{Z}^{M_Z}\} \\ \vdots & \ddots & \vdots \\ \Pr \{X_i \in \mathcal{X}^{M_X} | D_i = 0, Z_i \in \mathcal{Z}^1\} & \cdots & \Pr \{X_i \in \mathcal{X}^{M_X} | D_i = 1, Z_i \in \mathcal{Z}^{M_Z}\} \end{pmatrix}$$

We can apply the Kleibergen-Paap rank test or the eigenvalue ratio estimator.

Kleibergen and Paap (2006); Ahn and Horenstein (2013)

## Derivation of feasible moment function

$$\begin{aligned}
 & F_{Y(1)-Y(0)}(\delta) \\
 &= \mathbf{E} [\Pr \{Y_i(1) \leq Y_i(0) + \delta | U_i\}] \\
 &= \sum_{k=1}^K p_U(k) \underbrace{\int_{\mathbb{R}} \int_{\mathbb{R}} \mathbf{1}\{y \leq y' + \delta\} f_{Y(1)|U}(y|u^k) f_{Y(0)|U}(y'|u^k) dy dy'}_{=\Pr\{Y_i(1) \leq Y_i(0) + \delta | U_i = u^k\}} \quad \because Y_i(1) \perp\!\!\!\perp Y_i(0) \mid U_i \\
 &= \sum_{k=1}^K p_U(k) \int_{\mathbb{R}} \int_{\mathbb{R}} \mathbf{1}\{y \leq y' + \delta\} \left( \sum_{j=1}^K \tilde{\lambda}_{jk,1} f_{Y|D=1,Z}(y|z^j) \right) \quad \because \text{multiplying } \tilde{\Lambda}_d \text{ to } \mathbf{H}_d = \Gamma_d \cdot \Lambda_d \\
 &\quad \left( \sum_{j'=1}^K \tilde{\lambda}_{j'k,0} f_{Y|D=0,Z}(y'|z^{j'}) \right) dy dy' \\
 &= \sum_{k=1}^K \sum_{j=1}^K \sum_{j'=1}^K p_U(u^k) \tilde{\lambda}_{jk,1} \tilde{\lambda}_{j'k,0} \int_{\mathbb{R}} \int_{\mathbb{R}} \mathbf{1}\{y \leq y' + \delta\} f_{Y|D=1,Z}(y|z^j) f_{Y|D=0,Z}(y'|z^{j'}) dy dy'
 \end{aligned}$$

## Principal component analysis vs. nonnegative matrix factorization

Principal component analysis:

- given a  $M \times K$  matrix  $\mathbf{H}$  and an integer  $R > 0$ , find a rank  $R$  matrix  $\tilde{\mathbf{H}}$  such that

$$\min \left\| \mathbf{H} - \tilde{\mathbf{H}} \right\|_F$$

Nonnegative matrix factorization:

- given a  $M \times K$  matrix  $\mathbf{H}$  and an integer  $R > 0$ , find rank  $R$  **nonnegative** matrices  $\Gamma, \Lambda$  such that

$$\min \left\| \mathbf{H} - \Gamma \cdot \Lambda \right\|_F$$

NMF adds one more constraint: the low-rank representation should factor into nonnegative matrices.

## Nonnegative matrix factorization

$\Gamma_d$  can be further decomposed into  $\Gamma_X$  and  $\Gamma_{Y(d)}$ , using  $Y_i(d) \perp\!\!\!\perp X_i \mid U_i$ .

The minimization problem

$$\min \|\mathbb{H}_d - \Gamma_d \cdot \Lambda_d\|_F$$

becomes a quadratic program with linear constraints,  
once we fix two out of the three matrices  $\Gamma_X, \Gamma_{Y(d)}, \Lambda_d$ .

Thus, find the (local) minima by iterating across three objects:

1. Given  $(\Gamma_0^{(s)}, \Gamma_1^{(s)})$ , update  $(\Lambda_0, \Lambda_1)$ .
2. Given  $(\Lambda_0^{(s+1)}, \Lambda_1^{(s+1)}, \Gamma_{Y(0)}^{(s)}, \Gamma_{Y(1)}^{(s)})$ , update  $\Gamma_X$ .
3. Given  $(\Lambda_0^{(s+1)}, \Lambda_1^{(s+1)}, \Gamma_X^{(s+1)})$ , update  $(\Gamma_{Y(0)}, \Gamma_{Y(1)})$ .
4. Iterate **1-3** until convergence.

In practice, use many initial values to find the global minimum.

[back](#)

## Implementation: comparison to diagonalization

1. For each  $d, y$ , construct

$$\mathbb{H}_d(y) \left( \sum_{y'} \mathbb{H}_d(y') \right)^{-1}$$

where  $\mathbb{H}_d(y)$  estimates  $\Pr\{Y_i = y, X_i = x | D_i = d, Z_i = z\}$  and

$\sum_{y'} \mathbb{H}_d(y')$  estimates  $\Pr\{X_i = x | D_i = d, Z_i = z\}$ .

2. Diagonalize  $\mathbb{H}_d(y) \left( \sum_{y'} \mathbb{H}_d(y') \right)^{-1}$  across  $d, y$  since

$$\mathbf{H}_d(y) \left( \sum_{y'} \mathbf{H}_d(y') \right)^{-1} = \Gamma_X \cdot \Delta_d(y) \cdot \left( \Gamma_X \right)^{-1}.$$

Sum-to-one will pin down eigenvectors, i.e.  $\Gamma_X$ .



## Asymptotic normality + Delta method

The first-step NMF can be thought of as a (loosely defined) GMM or MLE estimator.

1. Nonnegativity constraints often bind.

The estimators may not be asymptotically normal.

e.g. Bonhomme et al. (2016) derives asymptotic normality, while not imposing nonnegativity.

2. DTE parameters are highly nonlinear in  $\Lambda_0$  and  $\Lambda_1$ .

Asymptotic normality using Delta method may converge slowly.

Bootstrapping for standard error  $\Rightarrow$  higher computation burden, given NMF.

3. Orthogonalization may help with discretization bias when  $U_i$  is continuous.

We would need  $\hat{\Lambda}_d$  and  $(\hat{\Lambda}_d)^{-1}$  to converge to true bivariate functions at  $n^{-\frac{1}{4}}$  rate.

Hence,  $\sqrt{n}$ -consistency and orthogonalization. [back](#)

## Additional moments in orthogonalization

$\Pr \{Y_i = y, X_i = x | U_i = u\} = \Pr \{Y_i = y | U_i = u\} \cdot \Pr \{X_i = x | U_i = u\}$  implies

$$\begin{aligned} & \frac{1}{2} \sum_{l=1}^K \frac{\tilde{\lambda}_{lk,d}}{p_{D,Z}(d,l)} \cdot \mathbf{E} \left[ \mathbf{1} \{Y_i = y, D_i = d, X_i = x, Z_i = z^l\} \right] \\ & + \frac{1}{2} \sum_{m=1}^K \frac{\tilde{\lambda}_{mk,d}}{p_{D,Z}(d,m)} \cdot \mathbf{E} \left[ \mathbf{1} \{Y_j = y, D_j = d, X_j = x, Z_i = z^m\} \right] \\ & - \frac{1}{2} \sum_{l=1}^K \sum_{m=1}^K \frac{\tilde{\lambda}_{lk,d} \tilde{\lambda}_{mk,d}}{p_{D,Z}(d,l) \cdot p_{D,Z}(d,m)} \mathbf{E} \left[ \mathbf{1} \{Y_i = y, D_i = d, Z_i = z^l, X_j = x, D_j = d, Z_j = z^m\} \right] \\ & - \frac{1}{2} \sum_{l=1}^K \sum_{m=1}^K \frac{\tilde{\lambda}_{lk,d} \tilde{\lambda}_{mk,d}}{p_{D,Z}(d,l) \cdot p_{D,Z}(d,m)} \mathbf{E} \left[ \mathbf{1} \{X_i = x, D_i = d, Z_i = z^l, Y_j = y, D_j = d, Z_j = z^m\} \right] = 0 \end{aligned}$$

with  $(Y_i, D_i, Z_i) \perp\!\!\!\perp (Y_j, D_j, Z_j)$ . [back](#)

## Falsification test

$Y_i(1) \perp\!\!\!\perp Y_i(0) \mid U_i$  from Assumption 2 is fundamentally untestable.

Instead, I test  $X_i \perp\!\!\!\perp D_i \mid U_i$  with estimators assuming  $Y_i(d) \perp\!\!\!\perp X_i \mid U_i$ .

“Can we construct a latent variable  $U_i$  that satisfies 1) conditional independence  $Y_i(d) \perp\!\!\!\perp X_i \mid U_i$  and 2) random treatment  $X_i \perp\!\!\!\perp D_i \mid U_i$ ?”

For this test, do not impose  $X_i \perp\!\!\!\perp D_i \mid U_i$  in the NMF.

In the short panel context,

- cannot test the conditional independence *across treatment regime*.
- can somewhat test the *intertemporal* conditional independence, given random treatment.

## Degrees of freedom

If testing  $X_i \perp\!\!\!\perp D_i \mid Z_i$  with

$$\hat{f}_{X|D=d,Z}(x|z) = \frac{\sum_{i=1}^n \mathbf{1}\{X_i = x, D_i = d, Z_i = z\}}{\sum_{i=1}^n \mathbf{1}\{D_i = d, Z_i = z\}}$$

$T_n \xrightarrow{d} \chi^2(K \cdot (M_X - 1))$  due to linear relationship  $\sum_{m=1}^{M_X} \hat{f}_{X|D=d,Z}(x^m|z) = 1$ .

$\hat{f}_{X|D=1,U}(x|u)$  and  $\hat{f}_{X|D=0,U}(x|u)$  uses additional moments  $\phi$  for orthogonalization.

Not necessarily  $\sum_{m=1}^{M_X} \hat{f}_{X|D=d,U}(x^m|u) = 1 \ \forall d, u$ . [back](#)

## Data generating process

The specifics of the DGPs are as follows:

- $n = 750, 2000$ .
- $D_i \perp\!\!\!\perp (Y_i(1), Y_i(0), X_i, Z_i, U_i)$  and  $\Pr\{D_i = 1\} = 0.5$ .
- $(p_U(1), p_U(2), p_U(3)) = (0.3, 0.3, 0.4)$ .
- $Y_i(d) \mid U_i = k \sim \mathcal{N}(\mu^k(d), \sigma^k(d)^2)$  and

$$\left(\mu^k(0), \sigma^k(0)\right) = \begin{cases} (-1, 1) & \text{if } k = 1 \\ (0, 1) & \text{if } k = 2 \\ (1, 1) & \text{if } k = 3 \end{cases} \quad \text{and} \quad \left(\mu^k(1), \sigma^k(1)\right) = \begin{cases} (1.5, 1.5) & \text{if } k = 1 \\ (2, 1) & \text{if } k = 2 \\ (2.5, 0.5) & \text{if } k = 3 \end{cases}.$$

- Since  $Y_i$  is continuous, a three-way partition is used:  $(-\infty, 0], (0, 2], (2, \infty)$ .

## Data generating process

Conditional distribution of  $Y_i(1) - Y_i(0)$  given  $U_i$ :

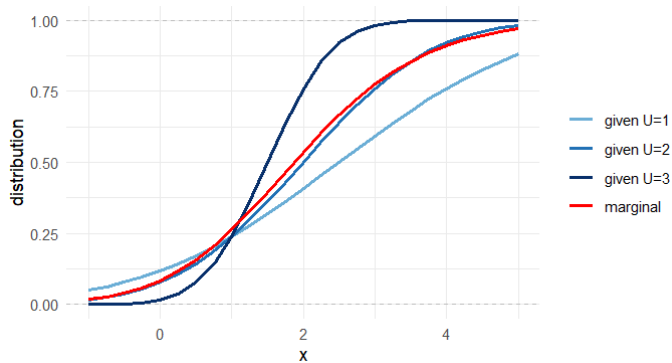


Figure 5: Marginal and conditional distributions of  $Y_i(1) - Y_i(0)$ .

## References I

- Ahn, Seung C and Alex R Horenstein**, "Eigenvalue ratio test for the number of factors," *Econometrica*, 2013, 81 (3), 1203–1227.
- Ai, Chunrong and Xiaohong Chen**, "Efficient estimation of models with conditional moment restrictions containing unknown functions," *Econometrica*, 2003, 71 (6), 1795–1843.
- Athey, Susan and Guido W Imbens**, "Identification and inference in nonlinear difference-in-differences models," *Econometrica*, 2006, 74 (2), 431–497.
- Attanasio, Orazio, Sarah Cattan, Emla Fitzsimons, Costas Meghir, and Marta Rubio-Codina**, "Estimating the production function for human capital: results from a randomized controlled trial in Colombia," *American Economic Review*, 2020, 110 (1), 48–85.
- Banerjee, Abhijit V, Shawn Cole, Esther Duflo, and Leigh Linden**, "Remedying education: Evidence from two randomized experiments in India," *The quarterly journal of economics*, 2007, 122 (3), 1235–1264.
- Bonhomme, Stéphane, Koen Jochmans, and Jean-Marc Robin**, "Estimating multivariate latent-structure models," *The Annals of Statistics*, 2016, pp. 540–563.
- Bonhomme, Stéphane, Thibaut Lamadon, and Elena Manresa**, "Discretizing unobserved heterogeneity," *Econometrica*, 2022, 90 (2), 625–643.
- Callaway, Brantly and Tong Li**, "Quantile treatment effects in difference in differences models with panel data," *Quantitative Economics*, 2019, 10 (4), 1579–1618.
- Carneiro, Pedro, Karsten T. Hansen, and James J. Heckman**, "2001 Lawrence R. Klein Lecture Estimating Distributions of Treatment Effects with an Application to the Returns to Schooling and Measurement of the Effects of Uncertainty on College Choice\*," *International Economic Review*, 2003, 44 (2), 361–422.

## References II

- Chen, Xiaohong and Xiaotong Shen**, “Sieve extremum estimates for weakly dependent data,” *Econometrica*, 1998, pp. 289–314.
- Chernozhukov, Victor and Christian Hansen**, “An IV model of quantile treatment effects,” *Econometrica*, 2005, 73 (1), 245–261.
- Chernozhukov, Victor and Christian Hansen**, “Instrumental quantile regression inference for structural and treatment effect models,” *Journal of Econometrics*, 2006, 132 (2), 491–525.
- Chernozhukov, Victor, Sokbae Lee, Adam M Rosen, and Liyang Sun**, “Policy Learning with Confidence,” *arXiv preprint arXiv:2502.10653*, 2025.
- Cunha, Flavio, James J Heckman, and Susanne M Schennach**, “Estimating the technology of cognitive and noncognitive skill formation,” *Econometrica*, 2010, 78 (3), 883–931.
- Deaner, Ben**, “Proxy controls and panel data,” 2023.
- Epstein, Larry G and Uzi Segal**, “Quadratic social welfare functions,” *Journal of Political Economy*, 1992, 100 (4), 691–712.
- Fan, Yanqin and Sang Soo Park**, “Sharp bounds on the distribution of treatment effects and their statistical inference,” *Econometric Theory*, 2010, 26 (3), 931–951.
- Fan, Yanqin, Robert Sherman, and Matthew Shum**, “Identifying treatment effects under data combination,” *Econometrica*, 2014, 82 (2), 811–822.
- Firpo, Sergio and Geert Ridder**, “Partial identification of the treatment effect distribution and its functionals,” *Journal of Econometrics*, 2019, 213 (1), 210–234.



## References III

- Frandsen, Brigham R and Lars J Lefgren**, "Partial identification of the distribution of treatment effects with an application to the Knowledge is Power Program (KIPP)," *Quantitative Economics*, 2021, 12 (1), 143–171.
- Heckman, James J and Edward Vytlacil**, "Structural equations, treatment effects, and econometric policy evaluation 1," *Econometrica*, 2005, 73 (3), 669–738.
- Heckman, James J, Jeffrey Smith, and Nancy Clements**, "Making the most out of programme evaluations and social experiments: Accounting for heterogeneity in programme impacts," *The Review of Economic Studies*, 1997, 64 (4), 487–535.
- Henry, Marc, Yuichi Kitamura, and Bernard Salanié**, "Partial identification of finite mixtures in econometric models," *Quantitative Economics*, 2014, 5 (1), 123–144.
- Hu, Yingyao**, "Identification and estimation of nonlinear models with misclassification error using instrumental variables: A general solution," *Journal of Econometrics*, 2008, 144 (1), 27–61.
- Hu, Yingyao and Matthew Shum**, "Nonparametric identification of dynamic models with unobserved state variables," *Journal of Econometrics*, 2012, 171 (1), 32–44.
- Hu, Yingyao and Susanne M Schennach**, "Instrumental variable treatment of nonclassical measurement error models," *Econometrica*, 2008, 76 (1), 195–216.
- Jones, Damon, David Molitor, and Julian Reif**, "What do workplace wellness programs do? Evidence from the Illinois workplace wellness study," *The Quarterly Journal of Economics*, 2019, 134 (4), 1747–1791.
- Kaji, Tetsuya and Jianfei Cao**, "Assessing Heterogeneity of Treatment Effects," 2023.

## References IV

- Kasahara, Hiroyuki and Katsumi Shimotsu**, “Nonparametric identification of finite mixture models of dynamic discrete choices,” *Econometrica*, 2009, 77 (1), 135–175.
- Kleibergen, Frank and Richard Paap**, “Generalized reduced rank tests using the singular value decomposition,” *Journal of econometrics*, 2006, 133 (1), 97–126.
- Levy, H and HM Markowitz**, “Approximating Expected Utility by a Function of Mean and Variance,” *The American Economic Review*, 1979, 69 (3), 308–317.
- Miao, Wang, Zhi Geng, and Eric J Tchetgen Tchetgen**, “Identifying causal effects with proxy variables of an unmeasured confounder,” *Biometrika*, 2018, 105 (4), 987–993.
- Mogstad, Magne, Andres Santos, and Alexander Torgovitsky**, “Using instrumental variables for inference about policy relevant treatment parameters,” *Econometrica*, 2018, 86 (5), 1589–1619.
- Muralidharan, Karthik, Abhijeet Singh, and Alejandro J Ganimian**, “Disrupting education? Experimental evidence on technology-aided instruction in India,” *American Economic Review*, 2019, 109 (4), 1426–1460.
- Nagasawa, Kenichi**, “Treatment effect estimation with noisy conditioning variables,” *arXiv preprint arXiv:1811.00667*, 2022.
- Noh, Sungho**, “Nonparametric identification and estimation of heterogeneous causal effects under conditional independence,” *Econometric Reviews*, 2023, 42 (3), 307–341.
- Reeve, Henry WJ, Timothy I Cannings, and Richard J Samworth**, “Optimal subgroup selection,” *The Annals of Statistics*, 2023, 51 (6), 2342–2365.
- Shen, Xiaotong**, “On methods of sieves and penalization,” *The Annals of Statistics*, 1997, 25 (6), 2555–2591.
- Wu, Ximing and Jeffrey M Perloff**, “Information-theoretic deconvolution approximation of treatment effect distribution,” *Available at SSRN 903982*, 2006.