

LEARNING WITH INSTRUMENTAL AND PROXY VARIABLES

Krikamol Muandet

Max Planck Institute for Intelligent Systems
Tübingen, Germany

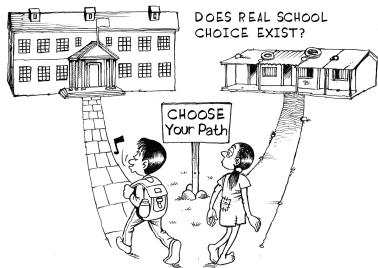
Joint work with

Jonas Kübler, Wittawat Jitkrittum, Arash Mehrjou, Si Kai Lee, Anant Raj, Rui Zhang,
Masaaki Imaizumi, Afsaneh Mastouri, Yuchen Zhu, Limor Gultchin, Anna Korba,
Bernhard Schölkopf, Ricardo Silva, Arthur Gretton, and Matt J. Kusner

Hi! PARIS Summer School 2022
July 7, 2022

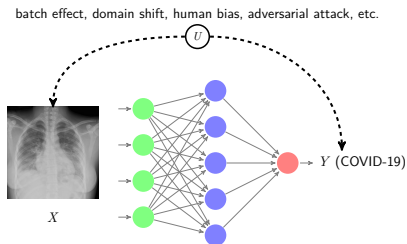
Unobserved Confounder / Spurious Correlation

Causal Inference



Education (X) \rightarrow Income (Y)

Machine Learning



$$Y = f(X) + \varepsilon_u$$

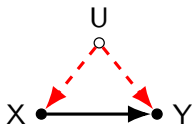
[UAI 2020; NeurIPS 2020; ICML 2021; Zhang et al. Under Review]

Spurious Correlation

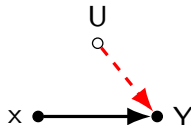
- ▶ An **interventional distribution** $P(Y \mid \text{do}(X = x))$.

Spurious Correlation

- ▶ An **interventional distribution** $P(Y \mid \text{do}(X = x))$.
- ▶ $P(Y \mid \text{do}(X = x)) \neq P(Y \mid X = x)$ **observational distribution**



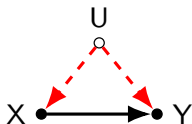
$$P(Y \mid X = x)$$



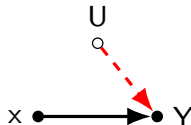
$$P(Y \mid \text{do}(X = x))$$

Spurious Correlation

- ▶ An **interventional distribution** $P(Y | \text{do}(X = x))$.
- ▶ $P(Y | \text{do}(X = x)) \neq P(Y | X = x)$ **observational distribution**

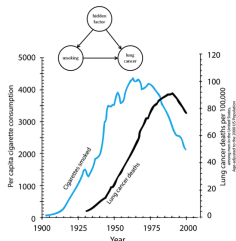


$$P(Y | X = x)$$



$$P(Y | \text{do}(X = x))$$

- ▶ The confounder U creates a **spurious correlation** between X and Y .



What if?



X = smoking (cigarettes smoked)

Y = lung cancer deaths

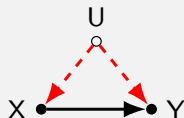
$P(Y | X = x)$ = observed deaths from data

$\text{do}(X = 0)$ = smoking banned

$P(Y | \text{do}(X = 0))$ = deaths after the ban

Additive Noise Model

Newey and Powell (2003); Hoyer et al. (2008); Peters et al. (2014)



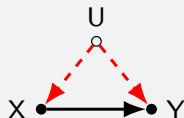
Structural Equations

$$Y \leftarrow f(X) + \varepsilon(U), \quad \mathbb{E}[\varepsilon] = 0$$

$$X \leftarrow g(U) + \nu, \quad \mathbb{E}[\nu] = 0$$

Additive Noise Model

Newey and Powell (2003); Hoyer et al. (2008); Peters et al. (2014)



Structural Equations

$$Y \leftarrow f(X) + \varepsilon(U), \quad \mathbb{E}[\varepsilon] = 0$$

$$X \leftarrow g(U) + \nu, \quad \mathbb{E}[\nu] = 0$$

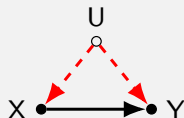
► **CI:** $\mathbb{E}[Y \mid \text{do}(X = x)] = f(x)$

ML: $\mathbb{E}[Y \mid X = x] = f(x)$

$$\mathbb{E}[Y \mid X = x] = f(x) + \mathbb{E}[\varepsilon \mid x]$$

Additive Noise Model

Newey and Powell (2003); Hoyer et al. (2008); Peters et al. (2014)



Structural Equations

$$Y \leftarrow f(X) + \varepsilon(U), \quad \mathbb{E}[\varepsilon] = 0$$

$$X \leftarrow g(U) + \nu, \quad \mathbb{E}[\nu] = 0$$

► **CI:** $\mathbb{E}[Y \mid \text{do}(X = x)] = f(x)$ **ML:** $\mathbb{E}[Y \mid X = x] = f(x)$

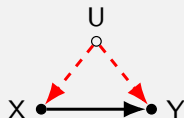
$$\mathbb{E}[Y \mid X = x] = f(x) + \mathbb{E}[\varepsilon \mid x]$$

► Estimate $\widehat{\mathbb{E}}[Y \mid X = x]$ from the sample $(x_1, y_1), \dots, (x_n, y_n) \sim P(X, Y)$

- A biased estimate of $\mathbb{E}[Y \mid \text{do}(X = x)]$.
- The estimate $\widehat{\mathbb{E}}[Y \mid X = x]$ can be **unstable** because of $\mathbb{E}[\varepsilon \mid x]$.

Additive Noise Model

Newey and Powell (2003); Hoyer et al. (2008); Peters et al. (2014)



Structural Equations

$$Y \leftarrow f(X) + \varepsilon(U), \quad \mathbb{E}[\varepsilon] = 0$$

$$X \leftarrow g(U) + \nu, \quad \mathbb{E}[\nu] = 0$$

► **CI:** $\mathbb{E}[Y \mid \text{do}(X = x)] = f(x)$ **ML:** $\mathbb{E}[Y \mid X = x] = f(x)$

$$\mathbb{E}[Y \mid X = x] = f(x) + \mathbb{E}[\varepsilon \mid x]$$

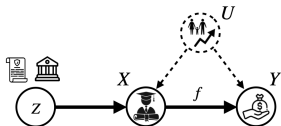
► Estimate $\widehat{\mathbb{E}}[Y \mid X = x]$ from the sample $(x_1, y_1), \dots, (x_n, y_n) \sim P(X, Y)$

► A biased estimate of $\mathbb{E}[Y \mid \text{do}(X = x)]$.

► The estimate $\widehat{\mathbb{E}}[Y \mid X = x]$ can be **unstable** because of $\mathbb{E}[\varepsilon \mid x]$.

► Impossible to recover $f(x)$ without further assumptions.

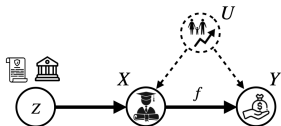
Instrumental Variables



Assumptions

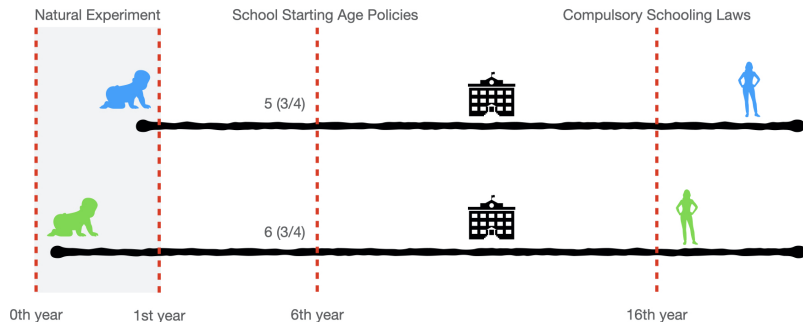
1. **Relevance:** $P(X|Z)$ is not constant in Z .
2. **Exclusion restriction:** Z affects Y only through X .
3. **Unconfoundedness:** Z is independent from U .

Instrumental Variables



Assumptions

1. **Relevance:** $P(X|Z)$ is not constant in Z .
2. **Exclusion restriction:** Z affects Y only through X .
3. **Unconfoundedness:** Z is independent from U .

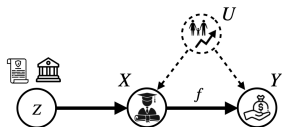


[Angrist and Krueger (1991); Angrist and Krueger (2001)]

Nobel Prize 2021

Nonparametric Instrumental Variable Regression

Newey and Powell (2003, *Econometrica*)



Assumptions

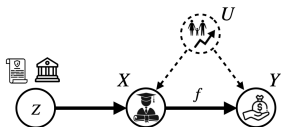
1. **Relevance:** $P(X|Z)$ is not constant in Z .
2. **Exclusion restriction:** Z affects Y only through X .
3. **Unconfoundedness:** Z is independent from U .

► Taking the conditional expectation wrt Z on both sides

$$\begin{aligned} Y &= f(X) + \varepsilon \\ \mathbb{E}[Y | Z] &= \mathbb{E}[f(X) | Z] + \underbrace{\mathbb{E}[\varepsilon | Z]}_{=\mathbb{E}[\varepsilon]=0} \end{aligned}$$

Nonparametric Instrumental Variable Regression

Newey and Powell (2003, *Econometrica*)



Assumptions

1. **Relevance:** $P(X|Z)$ is not constant in Z .
2. **Exclusion restriction:** Z affects Y only through X .
3. **Unconfoundedness:** Z is independent from U .

► Taking the conditional expectation wrt Z on both sides

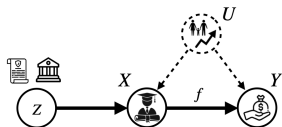
$$\begin{aligned} Y &= f(X) + \varepsilon \\ \mathbb{E}[Y | Z] &= \mathbb{E}[f(X) | Z] + \underbrace{\mathbb{E}[\varepsilon | Z]}_{=\mathbb{E}[\varepsilon]=0} \end{aligned}$$

► A **Fredholm integral equation** of the first kind

$$\mathbb{E}[Y | Z] = \int f(x) dP(x | Z)$$

Nonparametric Instrumental Variable Regression

Newey and Powell (2003, *Econometrica*)



Assumptions

1. **Relevance:** $P(X|Z)$ is not constant in Z .
2. **Exclusion restriction:** Z affects Y only through X .
3. **Unconfoundedness:** Z is independent from U .

- Taking the conditional expectation wrt Z on both sides

$$\begin{aligned} Y &= f(X) + \varepsilon \\ \mathbb{E}[Y | Z] &= \mathbb{E}[f(X) | Z] + \underbrace{\mathbb{E}[\varepsilon | Z]}_{=\mathbb{E}[\varepsilon]=0} \end{aligned}$$

- A **Fredholm integral equation** of the first kind


$$\mathbb{E}[Y | Z] = \int f(x) dP(x | Z)$$

- **Completeness condition** [D'Haultfoeuille (2011, *Econ Theory*)]

For all measurable functions g , $\mathbb{E}[g(X) | Z] = 0$ a.s. $\Rightarrow g(X) = 0$ a.s.

Fredholm Integral Equation: Two-Stage Estimation

Stage 1: Reduced Form


$$\mathbb{E}[Y | Z] = \int f(x) dP(x | Z)$$

Stage 2: Risk Minimization

Population risk

$$\min_{f \in \mathcal{F}} \mathbb{E}_Z \left[(\mathbb{E}[Y | Z] - \mathbb{E}_{X|Z}[f(X)])^2 \right]$$

Empirical risk: $(x_i, y_i, z_i) \sim P(X, Y, Z)$

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \left[(\hat{\mathbb{E}}[Y | z_i] - \hat{\mathbb{E}}_{X|z_i}[f(X)])^2 \right]$$

Fredholm Integral Equation: Two-Stage Estimation

Stage 1: Reduced Form

$$\mathbb{E}[Y | Z] = \int f(x) dP(x | Z)$$

Stage 2: Risk Minimization

Population risk

$$\min_{f \in \mathcal{F}} \mathbb{E}_Z \left[(\mathbb{E}[Y | Z] - \mathbb{E}_{X|Z}[f(X)])^2 \right]$$

Empirical risk: $(x_i, y_i, z_i) \sim P(X, Y, Z)$

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \left[(\hat{\mathbb{E}}[Y | z_i] - \hat{\mathbb{E}}_{X|z_i}[f(X)])^2 \right]$$

Previous work

- ▶ **2SLS** [Angrist and Imbens (1996)]
- ▶ **SieveIV** [Newey and Powell (2003)]
- ▶ **DeepIV** [Hartford (2017, ICML)]
- ▶ **KernelIV** [Singh (2019, NeurIPS)]

Fredholm Integral Equation: Two-Stage Estimation

Stage 1: Reduced Form

$$\mathbb{E}[Y | Z] = \int f(x) dP(x | Z)$$

Stage 2: Risk Minimization

Population risk

$$\min_{f \in \mathcal{F}} \mathbb{E}_Z [(\mathbb{E}[Y | Z] - \mathbb{E}_{X|Z}[f(X)])^2]$$

Empirical risk: $(x_i, y_i, z_i) \sim P(X, Y, Z)$

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n [(\hat{\mathbb{E}}[Y | z_i] - \hat{\mathbb{E}}_{X|z_i}[f(X)])^2]$$

Previous work

- ▶ **2SLS** [Angrist and Imbens (1996)]
- ▶ **SieveIV** [Newey and Powell (2003)]
- ▶ **DeepIV** [Hartford (2017, ICML)]
- ▶ **KernelIV** [Singh (2019, NeurIPS)]

Challenges

- ▶ Require **data splitting** for two-stage estimation.
- ▶ Estimate $P(X | Z)$ using **one observation** from $P(X | Z = z)$.
- ▶ The first stage is known as “**forbidden regression**” in econometrics.
- ▶ It violates **Vapnik's principle** [Vapnik 1998].

DualIV: From Two-Stage Estimation to Two-Player Game

Muandet, Mehrjou, Lee, Raj. (2020, **NeurIPS**); Liao et al (2020, **NeurIPS**)

$$\begin{aligned}\min_{f \in \mathcal{F}} R(f) &= \min_{f \in \mathcal{F}} \mathbb{E}_Z[(\mathbb{E}[Y | Z] - \mathbb{E}_{X|Z}[f(X)])^2] \\ &\stackrel{(a)}{=} \min_{f \in \mathcal{F}} \mathbb{E}_Z[\max_{u \in \mathbb{R}} \{\mathbb{E}_{X|Z}[f(X)]u - \mathbb{E}[Y|Z]u - \frac{1}{2}u^2\}] \\ &\stackrel{(b)}{=} \min_{f \in \mathcal{F}} \max_{u \in \mathcal{U}} \mathbb{E}_Z[\mathbb{E}_{X|Z}[f(X)]u(Z) - \mathbb{E}[Y|Z]u(Z) - \frac{1}{2}u^2(Z)] \\ &= \min_{f \in \mathcal{F}} \max_{u \in \mathcal{U}} \mathbb{E}_Z[(\mathbb{E}_{X|Z}[f(X)] - \mathbb{E}[Y|Z])u(Z) - \frac{1}{2}u^2(Z)] \\ &= \min_{f \in \mathcal{F}} \max_{u \in \mathcal{U}} \mathbb{E}_{XYZ}[(f(X) - Y)u(Z) - \frac{1}{2}u^2(Z)]\end{aligned}$$

(a) **Fenchel duality:** Let $\ell_y(\cdot) = \ell(y, \cdot)$ be a convex loss. For $\ell(y, y') = (y - y')^2$,

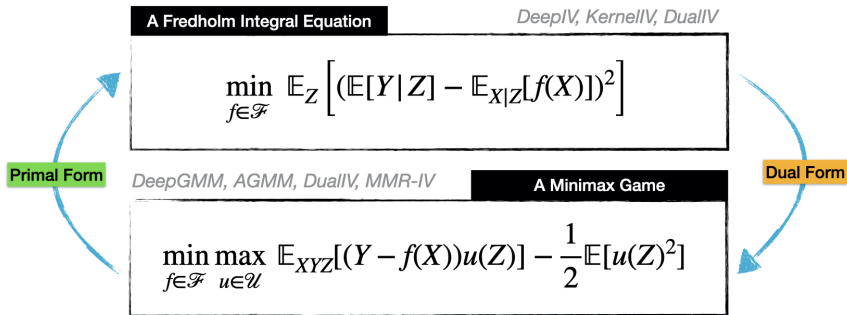
$$\ell_y(v) = \max_u \{uv - \ell_y^*(u)\}, \quad \ell_y^*(u) = uy + \frac{1}{2}u^2.$$

(b) **Interchangeability:** $\mathbb{E}_\omega[\max_{u \in \mathbb{R}} f(u, \omega)] = \max_{u(\cdot) \in \mathcal{U}} \mathbb{E}_\omega[f(u(\omega), \omega)]$.

Dai et al. (2017; Lemma 1), Rockafellar and Wets (1998; Ch. 14), and Shapiro et al. (2014; Ch. 7)

DualIV: From Two-Stage Estimation to Two-Player Game

Muandet, Mehrjou, Lee, Raj. (2020, **NeurIPS**); Liao et al (2020, **NeurIPS**)



[Liao et al (2020, **NeurIPS**); Bennett et al. (2019, **NeurIPS**); Muandet et al. (2020, **UAI**);
Dikkala et al. (2020, **NeurIPS**); Zhang et al. (2022, **Under Review**)]

DualIV: From Two-Stage Estimation to Two-Player Game

Muandet, Mehrjou, Lee, Raj. (2020, **NeurIPS**); Liao et al (2020, **NeurIPS**)

Data: $(x_1, y_1, z_1), \dots, (x_n, y_n, z_n) \sim P(X, Y, Z)$

► **Adversarial learning** ($\mathcal{F} = \{f_\psi : \psi \in \mathbb{R}^q\}$ and $\mathcal{U} = \{u_\theta : \theta \in \mathbb{R}^p\}$)

$$\min_{\psi \in \mathbb{R}^q} \max_{\theta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n [(y_i - f_\psi(x_i))u_\theta(z_i)] - \frac{\gamma}{2} \|u_\theta\|^2 + \frac{\lambda}{2} \|f_\psi\|^2$$

[Bennett et al. (2019, **NeurIPS**); Muandet et al. (2020, **UAI**); Dikkala et al. (2020, **NeurIPS**); Zhang et al. (2022, **UR**)]

DualIV: From Two-Stage Estimation to Two-Player Game

Muandet, Mehrjou, Lee, Raj. (2020, **NeurIPS**); Liao et al (2020, **NeurIPS**)

Data: $(x_1, y_1, z_1), \dots, (x_n, y_n, z_n) \sim P(X, Y, Z)$

- **Adversarial learning** ($\mathcal{F} = \{f_\psi : \psi \in \mathbb{R}^q\}$ and $\mathcal{U} = \{u_\theta : \theta \in \mathbb{R}^p\}$)

$$\min_{\psi \in \mathbb{R}^q} \max_{\theta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n [(y_i - f_\psi(x_i)) u_\theta(z_i)] - \frac{\gamma}{2} \|u_\theta\|^2 + \frac{\lambda}{2} \|f_\psi\|^2$$

- **M-estimators** (A unit ball within an RKHS, i.e., $\mathcal{U} = \{u : \mathcal{H}_{k_Z}, \|u\| \leq 1\}$)

$$\min_{f \in \mathcal{F}} R_V(f) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n (y_i - f(x_i)) k_Z(z_i, z_j) (y_j - f(x_j)) + \frac{\lambda}{2} \|f\|^2$$

[Bennett et al. (2019, **NeurIPS**); Muandet et al. (2020, **UAI**); Dikkala et al. (2020, **NeurIPS**); Zhang et al. (2022, **UR**)]

DualIV: From Two-Stage Estimation to Two-Player Game

Muandet, Mehrjou, Lee, Raj. (2020, **NeurIPS**); Liao et al (2020, **NeurIPS**)

Data: $(x_1, y_1, z_1), \dots, (x_n, y_n, z_n) \sim P(X, Y, Z)$

- **Adversarial learning** ($\mathcal{F} = \{f_\psi : \psi \in \mathbb{R}^q\}$ and $\mathcal{U} = \{u_\theta : \theta \in \mathbb{R}^p\}$)

$$\min_{\psi \in \mathbb{R}^q} \max_{\theta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n [(y_i - f_\psi(x_i))u_\theta(z_i)] - \frac{\gamma}{2} \|u_\theta\|^2 + \frac{\lambda}{2} \|f_\psi\|^2$$

- **M-estimators** (A unit ball within an RKHS, i.e., $\mathcal{U} = \{u : \mathcal{H}_{k_Z}, \|u\| \leq 1\}$)

$$\min_{f \in \mathcal{F}} R_V(f) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n (y_i - f(x_i))k_Z(z_i, z_j)(y_j - f(x_j)) + \frac{\lambda}{2} \|f\|^2$$

- **Kernel machines** ($\mathcal{F} \in \mathcal{H}_{k_X}$ and $\mathcal{U} = \{u : \mathcal{H}_{k_Z}, \|u\| \leq 1\}$)

$$f(x) = \sum_{i=1}^n \alpha_i k_X(x_i, x), \quad \alpha = (\mathbf{K}\mathbf{L}\mathbf{K} + \lambda\mathbf{K})^{-1}\mathbf{K}\mathbf{L}\mathbf{y}$$

where $\mathbf{K}_{ij} = k_X(x_i, x_j)$, $\mathbf{L}_{ij} = k_Z(z_i, z_j)/n^2$, and $\mathbf{y} = (y_1, \dots, y_n)^\top$.

[Bennett et al. (2019, **NeurIPS**); Muandet et al. (2020, **UAI**); Dikkala et al. (2020, **NeurIPS**); Zhang et al. (2022, **UR**)

DualIV: From Two-Stage Estimation to Two-Player Game

Muandet, Mehrjou, Lee, Raj. (2020, *NeurIPS*)

Demand design: $Y = f(X) + \varepsilon$ where $X = (P, T, S)$ and $Z = (C, T, S)$

- ▶ Y : sale, P : price, (T, S) : time of year and customer sentiment
- ▶ Sale Y and price P are confounded by supply-side market forces
- ▶ C : supply cost shifter (instrument)

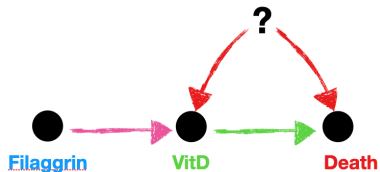
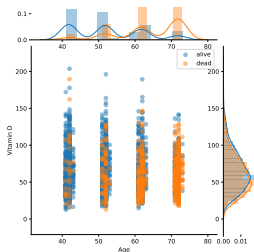
Table 1: Comparisons of IV regression methods in small (top) and medium (bottom) sample size regimes. We report the \log_{10} mean squared error (MSE) and its standard deviations over 20 trials.

	$n = 50$	Log_{10} Mean Squared Error (MSE)				
		$\rho = 0.1$	$\rho = 0.25$	$\rho = 0.5$	$\rho = 0.75$	$\rho = 0.9$
Two-Stage Estimation	2SLS	5.814 \pm 1.214	6.013 \pm 0.827	5.895 \pm 0.718	5.625 \pm 1.182	5.308 \pm 1.031
	DeepIV	5.127 \pm 0.043	5.131 \pm 0.031	5.133 \pm 0.072	5.130 \pm 0.124	5.127 \pm 0.061
	KernelIV	4.481 \pm 0.134	4.460 \pm 0.095	4.438 \pm 0.132	4.433 \pm 0.100	4.462 \pm 0.114
	DeepGMM	3.848 \pm 1.096	2.899 \pm 1.638	3.952 \pm 0.900	4.148 \pm 0.556	3.738 \pm 0.587
	DualIV	4.257 \pm 0.108	4.210 \pm 0.126	4.285 \pm 0.170	4.286 \pm 0.126	4.232 \pm 0.152
	$n = 1000$					
	2SLS	8.236 \pm 0.117	7.242 \pm 1.232	8.290 \pm 1.132	8.371 \pm 0.865	8.544 \pm 1.109
	DeepIV	4.613 \pm 0.052	4.618 \pm 0.048	4.614 \pm 0.068	4.701 \pm 0.040	4.731 \pm 0.032
	KernelIV	4.189 \pm 0.046	4.209 \pm 0.040	4.199 \pm 0.043	4.195 \pm 0.045	4.194 \pm 0.055
	DeepGMM	4.090 \pm 0.691	3.953 \pm 1.076	4.392 \pm 0.561	4.272 \pm 0.595	4.415 \pm 0.522
	DualIV	4.143 \pm 0.117	4.221 \pm 0.185	4.104 \pm 0.102	4.142 \pm 0.105	4.127 \pm 0.106

→ The strength of unobserved confounder →

Improvement over two-stage methods: $\sim 20\%$ ($n = 50$) and $\sim 9\%$ ($n = 1000$)

Vitamin D Data

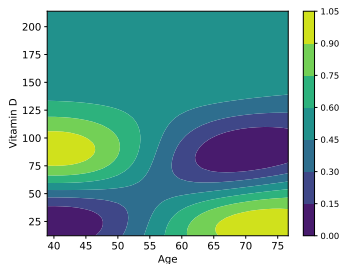


- ▶ A 10-year study on 2571 individuals aged 40–71.
- ▶ There are 4 variables:
 1. **Age** (at baseline)
 2. **Filaggrin** (binary indicator of filaggrin mutations)
 3. **VitD** (vitamin D level at baseline)
 4. **Death** (binary indicator of death during study)
- ▶ The goal is to evaluate the potential effect of VitD on Death.

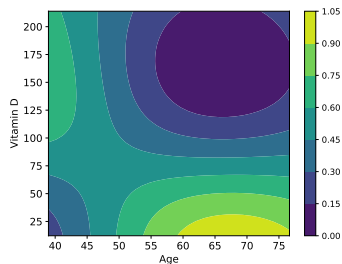
$$X = \text{VitD}, \quad Y = \text{Death}, \quad Z = \text{Filaggrin}$$

Vitamin D Data

$$\text{Death} = f(\text{VitD}, \text{Age})$$



(a) Kernel Ridge Regression (IV: None)

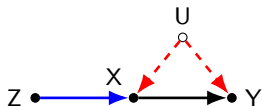


(b) Our Method (IV: Filaggrin Mutation)

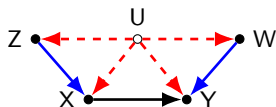
Proximal Causal Learning with Kernels

Mastouri, Zhu, Gultchin, Korba, Silva, Kusner, Gretton, Muandet (2021, **ICML**)

Instrumental Variable



Proxy Variables



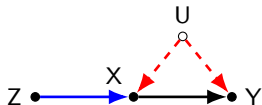
- Treatment-inducing proxy Z and outcome-inducing proxy W

[Miao et al. (2018, **Biometrika**)]

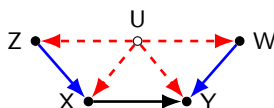
Proximal Causal Learning with Kernels

Mastouri, Zhu, Gultchin, Korba, Silva, Kusner, Gretton, Muandet (2021, **ICML**)

Instrumental Variable



Proxy Variables



- ▶ Treatment-inducing proxy Z and outcome-inducing proxy W
- ▶ The **causal effect** estimation

$$\mathbb{E}[Y \mid \text{do}(X = x)] = \mathbb{E}_W[h(x, W)]$$

where h is a **confounding bridge** function satisfying an integral equation

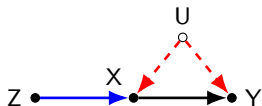
$$\mathbb{E}[Y \mid x, z] = \int h(x, w) dP(w \mid x, z)$$

[Miao et al. (2018, **Biometrika**)]

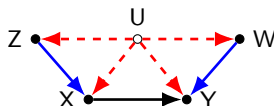
Proximal Causal Learning with Kernels

Mastouri, Zhu, Gultchin, Korba, Silva, Kusner, Gretton, Muandet (2021, **ICML**)

Instrumental Variable



Proxy Variables



- ▶ Treatment-inducing proxy Z and outcome-inducing proxy W
- ▶ The **causal effect** estimation

$$\mathbb{E}[Y \mid \text{do}(X = x)] = \mathbb{E}_W[h(x, W)]$$

where h is a **confounding bridge** function satisfying an integral equation

$$\mathbb{E}[Y \mid x, z] = \int h(x, w) dP(w \mid x, z)$$

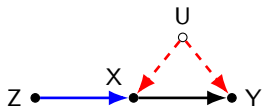
- ▶ The impact of **legalized abortion** (X) on **crime** (Y) [Donohue and Levitt (2001)]
(Supreme Court's 1973 decision in *Roe v. Wade*)

[Miao et al. (2018, **Biometrika**)]

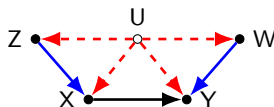
Proximal Causal Learning with Kernels

Mastouri, Zhu, Gultchin, Korba, Silva, Kusner, Gretton, Muandet (2021, ICML)

Instrumental Variable



Proxy Variables



- ▶ Treatment-inducing proxy Z and outcome-inducing proxy W
- ▶ The **causal effect** estimation

$$\mathbb{E}[Y \mid \text{do}(X = x)] = \mathbb{E}_W[h(x, W)]$$

where h is a **confounding bridge** function satisfying an integral equation

$$\mathbb{E}[Y \mid x, z] = \int h(x, w) dP(w \mid x, z)$$

- ▶ The impact of **legalized abortion** (X) on **crime** (Y) [Donohue and Levitt (2001)]
(Supreme Court's 1973 decision in *Roe v. Wade*)
- ▶ The impact of **grade retention** (X) on **cognitive outcome** (Y) [Deaner (2018)]

[Miao et al. (2018, Biometrika)]

Fantastic “Instruments” and Where to Find Them

Mendelian Randomization

[Adam (2019, *Nature*)]



Issue: weak instruments

Fantastic “Instruments” and Where to Find Them

Mendelian Randomization

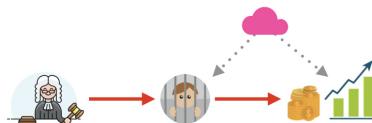
[Adam (2019, **Nature**)]



Issue: weak instruments

Judge Leniency Design

[Kling (2006, **AER**)]



Issue: algorithmic decision making

Fantastic “Instruments” and Where to Find Them

Mendelian Randomization

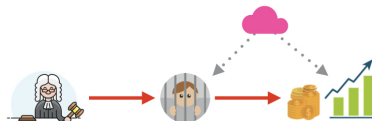
[Adam (2019, **Nature**)]



Issue: weak instruments

Judge Leniency Design

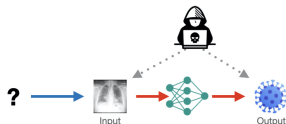
[Kling (2006, **AER**)]



Issue: algorithmic decision making

Adversarial Examples

[N.A.]



Issue: prone to adversarial attack

Fantastic “Instruments” and Where to Find Them

Mendelian Randomization

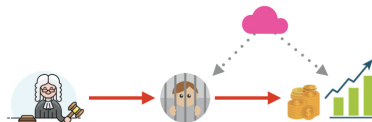
[Adam (2019, **Nature**)]



Issue: weak instruments

Judge Leniency Design

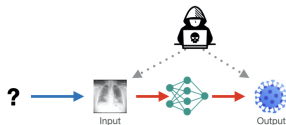
[Kling (2006, **AER**)]



Issue: algorithmic decision making

Adversarial Examples

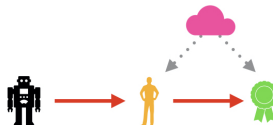
[N.A.]



Issue: prone to adversarial attack

Algorithmic Instruments

[Ngo et al. (2021, **ICML**)]



Issue: non-compliance