Causal Effect Estimation with Context and Confounders

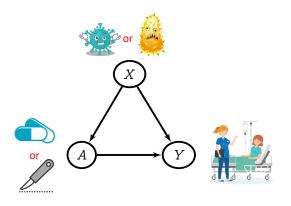
Arthur Gretton

Gatsby Computational Neuroscience Unit, UCL
Google Deepmind

Causality-XAI Winter School, 2025

Observation vs intervention

Conditioning from observation: $\mathbb{E}[Y|A=a] = \sum_{x} \mathbb{E}[Y|a,x] p(x|a)$

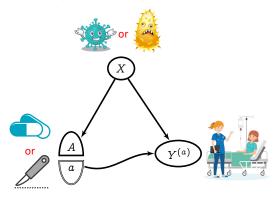


From our *observations* of historical hospital data:

- P(Y = cured|A = pills) = 0.85
- P(Y = cured|A = surgery) = 0.72

Observation vs intervention

Average causal effect (intervention): $\mathbb{E}[Y^{(a)}] = \sum_{x} \mathbb{E}[Y|a,x]p(x)$

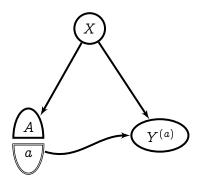


From our *intervention* (making all patients take a treatment):

- $P(Y^{(pills)} = cured) = 0.64$
- $P(Y^{(\text{surgery})} = \text{cured}) = 0.75$

Richardson, Robins (2013), Single World Intervention Graphs (SWIGs): A Unification of the Counterfactual and Graphical Approaches to Causality

Questions we will solve



Outline

Causal effect estimation, observed covariates:

 Average treatment effect (ATE), conditional average treatment effect (CATE)

Causal effect estimation, hidden covariates:

■ ... instrumental variables, proxy variables

What's new? What is it good for?

- Treatment A, covariates X, etc can be multivariate, complicated...
- ...by using kernel or adaptive neural net feature representations

Model assumption: linear functions of features

All learned functions will take the form:

$$oldsymbol{\gamma}(x) = oldsymbol{\gamma}^ op arphi(x) = \left_{\mathcal{H}}$$

Model assumption: linear functions of features

All learned functions will take the form:

$$oldsymbol{\gamma}(x) = oldsymbol{\gamma}^ op arphi(x) = \left_{\mathcal{H}}$$

Option 1: Finite dictionaries of learned neural net features $\varphi_{\theta}(x)$ (linear final layer γ)

Xu, G., A Neural mean embedding approach for back-door and front-door adjustment. (ICLR 23)

Xu, Chen, Srinivasan, de Freitas, Doucet, G. Learning Deep Features in Instrumental Variable Regression. (ICLR 21)

Option 2: Infinite dictionaries of fixed kernel features:

$$\left\langle arphi(x_i),arphi(x)
ight
angle_{\mathcal{H}}=k(x_i,x)$$

Kernel is feature dot product.

Singh, Xu, G. Kernel Methods for Causal Functions: Dose, Heterogeneous, and Incremental Response Curves. (Biometrika, 2023)

Singh, Sahani, G. Kernel Instrumental Variable Regression. (NeurIPS 19)

Model fitting: neural ridge regression

Learn $\gamma_0(x) := \mathbb{E}[Y|X=x]$ from features $\varphi_{\theta}(x_i)$ with outcomes y_i :

$$\hat{m{\gamma}} = rg \min_{m{\gamma} \in \mathcal{H}} \left(\sum_{i=1}^n \left(y_i - m{\gamma}^ op m{arphi}_ heta(m{x}_i)
ight)^2 + \lambda \|m{\gamma}\|^2
ight)$$
 (1

Model fitting: neural ridge regression

Learn $\gamma_0(x) := \mathbb{E}[Y|X=x]$ from features $\varphi_{\theta}(x_i)$ with outcomes y_i :

$$\hat{\gamma} = \arg\min_{\gamma \in \mathcal{H}} \left(\sum_{i=1}^{n} \left(y_i - \gamma^{\top} \varphi_{\theta}(x_i) \right)^2 + \lambda \|\gamma\|^2 \right)$$
 (1)

Solution for linear final layer γ :

$$egin{aligned} \hat{\gamma} &= C_{YX}^{(heta)} (\, C_{XX}^{(heta)} + \lambda)^{-1} \ C_{YX}^{(heta)} &= rac{1}{n} \sum_{i=1}^n [y_i \ arphi_{ heta}(x_i)^ op] \ C_{XX}^{(heta)} &= rac{1}{n} \sum_{i=1}^n [arphi_{ heta}(x_i) \ arphi_{ heta}(x_i) \ arphi_{ heta}(x_i)^ op] \end{aligned}$$

Model fitting: neural ridge regression

Learn $\gamma_0(x) := \mathbb{E}[Y|X=x]$ from features $\varphi_{\theta}(x_i)$ with outcomes y_i :

$$\hat{\gamma} = \arg\min_{\gamma \in \mathcal{H}} \left(\sum_{i=1}^{n} \left(y_i - \gamma^{\top} \varphi_{\theta}(x_i) \right)^2 + \lambda \|\gamma\|^2 \right)$$
 (1)

Solution for linear final layer γ :

$$egin{aligned} \hat{\gamma} &= C_{YX}^{(heta)} (\, C_{XX}^{(heta)} + \lambda)^{-1} \ C_{YX}^{(heta)} &= rac{1}{n} \sum_{i=1}^n [y_i \ arphi_{ heta}(x_i)^ op] \ C_{XX}^{(heta)} &= rac{1}{n} \sum_{i=1}^n [arphi_{ heta}(x_i) \ arphi_{ heta}(x_i) \ arphi_{ heta}(x_i)^ op] \end{aligned}$$

How to solve for θ :

Substitute $\hat{\gamma}$ into (1), backprop through Cholesky for θ .

More details: Galashov, Da Costa, Xu, Hennig, G, Closed-Form Last Layer Optimization (2025, arxiv:2510.04606)

Model fitting: kernel ridge regression

Learn $\gamma_0(x) := \mathbb{E}[Y|X=x]$ from features $\varphi(x_i)$ with outcomes y_i :

$$\hat{\gamma} = \arg\min_{\gamma \in \mathcal{H}} \left(\sum_{i=1}^{n} (y_i - \langle \gamma, \varphi(x_i) \rangle_{\mathcal{H}})^2 + \lambda \|\gamma\|_{\mathcal{H}}^2 \right).$$

Model fitting: kernel ridge regression

Learn $\gamma_0(x) := \mathbb{E}[Y|X=x]$ from features $\varphi(x_i)$ with outcomes y_i :

$$\hat{\gamma} = \arg\min_{\gamma \in \mathcal{H}} \left(\sum_{i=1}^{n} (y_i - \langle \gamma, \varphi(x_i) \rangle_{\mathcal{H}})^2 + \lambda \|\gamma\|_{\mathcal{H}}^2 \right).$$

Infinite dimensional solution at x:

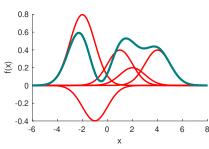
$$egin{aligned} \hat{\gamma}(x) &= C_{YX}(C_{XX} + \lambda)^{-1} arphi(x) \ C_{YX} &= rac{1}{n} \sum_{i=1}^n [y_i \ arphi(x_i)^ op] \ C_{XX} &= rac{1}{n} \sum_{i=1}^n [arphi(x_i) \ arphi(x_i)^ op] \end{aligned}$$

Model fitting: kernel ridge regression

Learn $\gamma_0(x) := \mathbb{E}[Y|X=x]$ from features $\varphi(x_i)$ with outcomes y_i :

$$\hat{\gamma} = \arg\min_{\gamma \in \mathcal{H}} \left(\sum_{i=1}^{n} (y_i - \langle \gamma, \varphi(x_i) \rangle_{\mathcal{H}})^2 + \lambda \|\gamma\|_{\mathcal{H}}^2 \right).$$

Kernel solution at x (as weighted sum of y) $\hat{\gamma}(x) = \sum_{i=1}^{n} y_i \beta_i(x)$ $\beta(x) = (K_{XX} + \lambda I)^{-1} k_{Xx}$ $(K_{XX})_{ij} = k(x_i, x_j) = \langle \varphi(x_i), \varphi(x_j) \rangle_{\mathcal{H}}$ $(k_{Xx})_i = k(x_i, x)$



Observed covariates: (conditional) ATE

Kernels (Biometrika 2023):







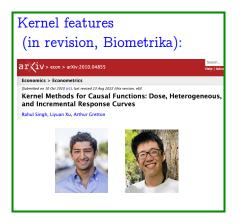
NN features (ICLR 2023):





Code for NN and kernel causal estimation with observed covariates: https://github.com/liyuan9988/DeepFrontBackDoor/

Observed covariates: (conditional) ATE



NN features (ICLR 2023):



Code for NN and kernel causal estimation with observed covariates: https://github.com/liyuan9988/DeepFrontBackDoor/ 9,

Average treatment effect

Potential outcome (intervention):

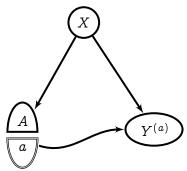
$$\mathbb{E}[\,Y^{(\,a)}] = \int \mathbb{E}[\,Y|\,a,x] \, dp(x)$$

(the average structural function; in epidemiology, for continuous a, the dose-response curve).

Assume: (1) Stable Unit Treatment Value Assumption (aka "no interference"), (2) Conditional exchangeability $Y^{(a)} \perp \!\!\!\perp A|X$. (3) Overlap.

Example: US job corps, training for disadvantaged youths:

- A: treatment (training hours)
- Y: outcome (percentage employment)
- X: covariates (age, education, marital status, ...)



Multiple inputs via products of kernels

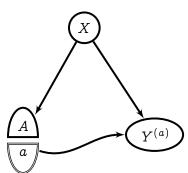
We may predict expected outcome from two inputs

$$\gamma_0(a,x) := \mathbb{E}[Y|a,x]$$

Assume we have:

- covariate features $\varphi(x)$ with kernel k(x, x')
- treatment features $\varphi(a)$ with kernel k(a, a')

(argument of kernel/feature map indicates feature space)



Multiple inputs via products of kernels

We may predict expected outcome from two inputs

$$\gamma_0(a,x) := \mathbb{E}[Y|a,x]$$

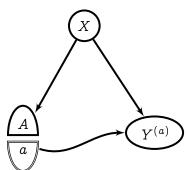
Assume we have:

- covariate features $\varphi(x)$ with kernel k(x, x')
- treatment features $\varphi(a)$ with kernel k(a, a')

(argument of kernel/feature map indicates feature space)

We use outer product of features (\Longrightarrow product of kernels):

$$\phi(x,a) = \varphi(a) \otimes \varphi(x)$$
 $\mathfrak{K}([a,x],[a',x']) = k(a,a')k(x,x')$



Multiple inputs via products of kernels

We may predict expected outcome from two inputs

$$\gamma_0(a,x) := \mathbb{E}[Y|a,x]$$

Assume we have:

- covariate features $\varphi(x)$ with kernel k(x, x')
- treatment features $\varphi(a)$ with kernel k(a, a')

(argument of kernel/feature map indicates feature space)

We use outer product of features (\implies product of kernels):

$$\phi(x,a)=arphi(a)\otimesarphi(x) \qquad \mathfrak{K}([a,x],[a',x'])=k(a,a')k(x,x')$$

a

Ridge regression solution:

$$\hat{\gamma}(x,a) = \sum_{i=1}^{n} y_i eta_i(a,x), \;\; eta(a,x) = \left[K_{AA} \odot K_{XX} + \lambda I
ight]^{-1} K_{Aa} \odot K_{XY}$$

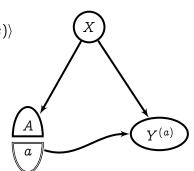
ATE (dose-response curve)

Well-specified setting:

$$\mathbb{E}[\,Y|\,a,x]=:\gamma_0(\,a,x)=\langle\gamma_0,arphi(\,a)\otimesarphi(\,x)
angle$$

ATE as feature space dot product:

$$egin{aligned} ext{ATE}(a) &= \mathbb{E}[\gamma_0(a,X)] \ &= \mathbb{E}\left[\langle \gamma_0, arphi(a) \otimes arphi(X)
angle
ight] \end{aligned}$$



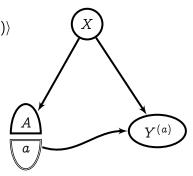
ATE (dose-response curve)

Well-specified setting:

$$\mathbb{E}[\,Y|\,a,x]=:\gamma_0(a,x)=\langle\gamma_0,arphi(a)\otimesarphi(x)
angle$$

ATE as feature space dot product:

$$egin{aligned} ext{ATE}(a) &= \mathbb{E}[\gamma_0(a,X)] \ &= \mathbb{E}\left[\langle \gamma_0, arphi(a) \otimes arphi(X)
angle
ight] \ &= \langle \gamma_0, arphi(a) \otimes \underbrace{\mu_X}_{\mathbb{E}[arphi(X)]}
angle \end{aligned}$$



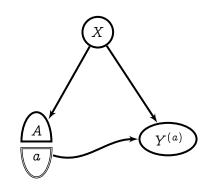
Feature map of probability P(X),

$$\mu_{X} = [\dots \mathbb{E}\left[\varphi_{i}(X)\right]\dots]$$

ATE: example

US job corps: training for disadvantaged youths:

- X: covariate/context (age, education, marital status, ...)
- A: treatment (training hours)
- Y: outcome (percent employment)



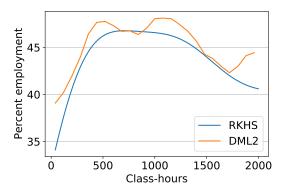
Empirical ATE:

$$egin{aligned} \widehat{ ext{ATE}}(a) &= \widehat{\mathbb{E}}\left[\left\langle \hat{\gamma}_0, arphi(X) \otimes arphi(a)
ight
angle
ight] \ &= rac{1}{n} \sum_{i=1}^n Y^ op (K_{AA} \odot K_{XX} + n\lambda I)^{-1} (K_{Aa} \odot K_{Xx_i}) \end{aligned}$$

Schochet, Burghardt, and McConnell (2008). Does Job Corps work? Impact findings from the national Job Corps study.

13/56
Singh, Xu. G (2023).

ATE: results



- First 12.5 weeks of classes confer employment gain: from 35% to 47%.
- [RKHS] is our $\widehat{ATE}(a)$.
- [DML2] Colangelo, Lee (2020), Double debiased machine learning nonparametric inference with continuous treatments.

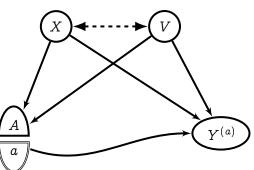
Singh, Xu, G (2023)

Well-specified setting:

$$egin{aligned} \mathbb{E}[\,Y|\,a,x,v] =: \gamma_0(\,a,x,v) \ &= \langle \gamma_0, arphi(\,a) \otimes arphi(x) \otimes arphi(v)
angle \,. \end{aligned}$$

Conditional ATE

$$=\mathbb{E}\left[\left.Y^{(a)}
ight| rac{oldsymbol{V}}{oldsymbol{v}}=oldsymbol{v}
ight]$$



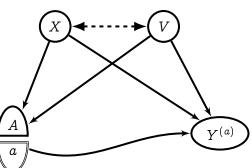
Well-specified setting:

$$\mathbb{E}[Y|a,x,v]=:\gamma_0(a,x,v) \ = \langle \gamma_0, arphi(a)\otimes arphi(x)\otimes arphi(v)
angle \,.$$

Conditional ATE

$$=\mathbb{E}\left[\left.Y^{\left(a
ight)}
ight|oldsymbol{V}=oldsymbol{v}
ight]$$

$$oxed{=\mathbb{E}\left[\left\langle \gamma_{0},arphi(a)\otimesarphi(X)\otimesarphi(extbf{ extit{V}})
ight
angle \leftert extbf{ extit{V}}= extbf{ extit{v}}
ight]}$$



Well-specified setting:

$$egin{aligned} \mathbb{E}[\,Y|\,a,x,v] &=: \gamma_0(a,x,v) \ &= \langle \gamma_0, arphi(a) \otimes arphi(x) \otimes arphi(v)
angle \,. \end{aligned}$$

Conditional ATE

CATE
$$(a, v)$$

$$= \mathbb{E}\left[Y^{(a)}|V = v\right]$$

$$=\mathbb{E}\left[\left\langle \gamma_{0},arphi(a)\otimesarphi(X)\otimesarphi(rac{V}{V})
ight
angle \left|rac{V}{V}
ight.$$

= ...?

How to take conditional expectation?

Density estimation for p(X|V=v)? Sample from p(X|V=v)?



Well-specified setting:

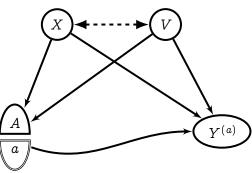
$$egin{aligned} \mathbb{E}[\,Y|\,a,x,v] &=: \gamma_0(\,a,x,v) \ &= \langle \gamma_0, arphi(\,a) \otimes arphi(x) \otimes arphi(v)
angle \,. \end{aligned}$$

Conditional ATE

$$\begin{aligned}
& \text{CATE}(a, v) \\
&= \mathbb{E}\left[Y^{(a)} | V = v\right] \\
&= \mathbb{E}\left[\langle \gamma_0, \varphi(a) \otimes \varphi(X) \otimes \varphi(V) \rangle | V = v\right] \\
&= \langle \gamma_0, \varphi(a) \otimes \mathbb{E}[\varphi(X) | V = v] \otimes \varphi(v) \rangle
\end{aligned}$$

 $\mu_{X|V=v}$

Learn conditional mean embedding: $\mu_{X|V=v} := \mathbb{E}_X \left[\varphi(X) \middle| V=v \right]$



Our goal: an operator $F_0: \mathcal{H}_{\mathcal{V}} \to \mathcal{H}_{\mathcal{X}}$ such that

$$F_0 \varphi(v) = \mu_{X|V=v}$$

Song, Huang, Smola, Fukumizu (2009). Hilbert space embeddings of conditional distributions with applications to dynamical systems.

Grunewalder, Lever, Baldassarre, Patterson, G, Pontil (2012). Conditional mean embeddings as regressors.

Grunewalder, G, Shawe-Taylor (2013) Smooth operators.

Li, Meunier, Mollenhauer, G (2022), Optimal Rates for Regularized Conditional Mean Embedding Learning

Our goal: an operator $F_0: \mathcal{H}_{\mathcal{V}} \to \mathcal{H}_{\mathcal{X}}$ such that

$$F_0 \varphi(v) = \mu_{X|V=v}$$

Assume

$$F_0 \in \overline{\operatorname{span}\left\{ \varphi(x) \otimes \varphi(v) \right\}} \iff F_0 \in \operatorname{HS}(\mathcal{H}_{\mathcal{V}}, \mathcal{H}_{\mathcal{X}})$$

Implied smoothness assumption:

$$\mathbb{E}[h(X)|V=v]\in\mathcal{H}_{\mathcal{V}}\quadorall h\in\mathcal{H}_{\mathcal{X}}$$

Song, Huang, Smola, Fukumizu (2009). Hilbert space embeddings of conditional distributions with applications to dynamical systems.

Grunewalder, Lever, Baldassarre, Patterson, G, Pontil (2012). Conditional mean embeddings as regressors.

Grunewalder, G, Shawe-Taylor (2013) Smooth operators.

Li, Meunier, Mollenhauer, G (2022), Optimal Rates for Regularized Conditional Mean Embedding Learning 16/56

Our goal: an operator $F_0: \mathcal{H}_{\mathcal{V}} \to \mathcal{H}_{\mathcal{X}}$ such that

$$F_0 \varphi(v) = \mu_{X|V=v}$$

Assume

$$F_0 \in \overline{\operatorname{span}\left\{arphi(x) \otimes arphi(v)
ight\}} \iff F_0 \in \operatorname{HS}(\mathcal{H}_\mathcal{V},\mathcal{H}_\mathcal{X})$$

Implied smoothness assumption:

$$\mathbb{E}[h(X)|\, rac{oldsymbol{V}}{oldsymbol{V}}] \in \mathcal{H}_{\mathcal{V}} \quad orall h \in \mathcal{H}_{\mathcal{X}}$$

A Smooth Operator

Song, Huang, Smola, Fukumizu (2009). Hilbert space embeddings of conditional distributions with applications to dynamical systems.

Grunewalder, Lever, Baldassarre, Patterson, G, Pontil (2012). Conditional mean embeddings as regressors.

Grunewalder, G, Shawe-Taylor (2013) Smooth operators.

Li, Meunier, Mollenhauer, G (2022), Optimal Rates for Regularized Conditional Mean Embedding Learning

Our goal: an operator $F_0: \mathcal{H}_{\mathcal{V}} \to \mathcal{H}_{\mathcal{X}}$ such that

$$F_0 \varphi(v) = \mu_{X|V=v}$$

Assume

$$F_0 \in \overline{\operatorname{span}\left\{ \varphi(x) \otimes \varphi(v) \right\}} \iff F_0 \in \operatorname{HS}(\mathcal{H}_{\mathcal{V}}, \mathcal{H}_{\mathcal{X}})$$

Implied smoothness assumption:

$$\mathbb{E}[h(X)| rac{oldsymbol{V}}{oldsymbol{V}} = rac{oldsymbol{v}}{oldsymbol{V}}] \in \mathcal{H}_{\mathcal{V}} \quad orall h \in \mathcal{H}_{\mathcal{X}}$$

Kernel ridge regression from $\varphi(v)$ to infinite features $\varphi(x)$:

$$\widehat{F} = \operatorname*{argmin}_{F \in HS} \sum_{\ell=1}^n \| arphi(x_\ell) - F arphi(v_\ell) \|_{\mathcal{H}_{\mathcal{X}}}^2 + \lambda_2 \| F \|_{HS}^2$$

Song, Huang, Smola, Fukumizu (2009). Hilbert space embeddings of conditional distributions with applications to dynamical systems.

Grunewalder, Lever, Baldassarre, Patterson, G, Pontil (2012). Conditional mean embeddings as regressors.

Grunewalder, G, Shawe-Taylor (2013) Smooth operators.

Li, Meunier, Mollenhauer, G (2022), Optimal Rates for Regularized Conditional Mean Embedding Learning 16/56

Our goal: an operator $F_0: \mathcal{H}_{\mathcal{V}} \to \mathcal{H}_{\mathcal{X}}$ such that

$$F_0\varphi(v)=\mu_{X|V=v}$$

Assume

$$F_0 \in \overline{\operatorname{span}\left\{ arphi(x) \otimes arphi(v)
ight\}} \iff F_0 \in \operatorname{HS}(\mathcal{H}_{\mathcal{V}},\mathcal{H}_{\mathcal{X}})$$

Implied smoothness assumption:

$$\mathbb{E}[h(X)| extbf{ extit{V}}= extbf{ extit{v}}]\in\mathcal{H}_{\mathcal{V}}\quadorall h\in\mathcal{H}_{\mathcal{X}}$$

Kernel ridge regression from $\varphi(v)$ to infinite features $\varphi(x)$:

$$\widehat{F} = \operatorname*{argmin}_{F \in HS} \sum_{\ell=1}^n \| arphi(x_\ell) - F arphi(v_\ell) \|_{\mathcal{H}_{\mathcal{X}}}^2 + \lambda_2 \| F \|_{HS}^2$$

Ridge regression solution:

$$egin{aligned} \mu_{X|V=oldsymbol{v}} := \mathbb{E}[arphi(X)|oldsymbol{V} = oldsymbol{v}] &pprox \widehat{F}arphi(oldsymbol{v}) = \sum_{\ell=1}^n arphi(x_\ell)eta_\ell(oldsymbol{v}) \ eta(oldsymbol{v}) = [K_{VV} + \lambda_2 I]^{-1} \, k_{Vv} \end{aligned}$$

16/56

Conditional ATE: example

US job corps:

- X: confounder/context (education, marital status, ...)
- A: treatment (training hours)
- Y: outcome (percent employed)
- *V*: age

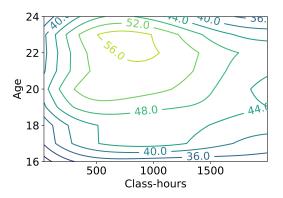
X V $Y^{(a)}$

Empirical CATE:

$$\widehat{ ext{CATE}}(a, extbf{v}) = \langle \hat{\gamma}_0, arphi(a) \otimes \underbrace{\widehat{F}arphi(extbf{v})}_{\widehat{\mathbb{E}}[arphi(extbf{x})]} \otimes arphi(extbf{v})
angle$$

(with consistency guarantees: see paper!)

Conditional ATE: results



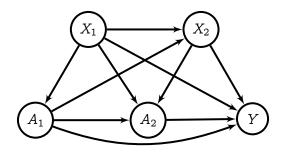
Average percentage employment $Y^{(a)}$ for class hours a, conditioned on age v. Given around 12-14 weeks of classes:

- 16 y/o: employment increases from 28% to at most 36%.
- 22 y/o: percent employment increases from 40% to 56%.

Singh, Xu, G (2023)

...dynamic treatment effect...

Dynamic treatment effect: sequence A_1 , A_2 of treatments.



- potential outcomes $Y^{(a_1)}$, $Y^{(a_2)}$, $Y^{(a_1,a_2)}$,
 counterfactuals $\mathbb{E}\left[Y^{(a'_1,a'_2)}|A_1=a_1,A_2=a_2\right]...$ (c.f. the Robins G-formula)

Singh, Xu, G. (Bernoulli 2025) Kernel Methods for Multistage Causal Inference: Mediation Analysis and **Dynamic Treatment Effects**

What if there are hidden confounders?

Ticket price A, seats sold Y.

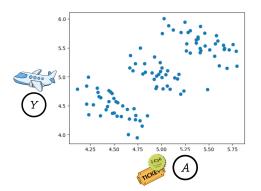


What is the effect on seats sold $Y^{(a)}$ of intervening on price a?

Ticket price A, seats sold Y.

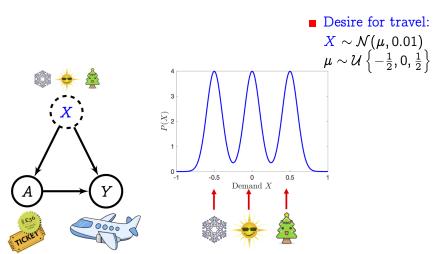


What is the effect on seats sold $Y^{(a)}$ of intervening on price a?

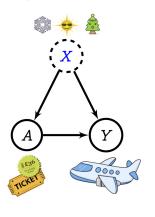


Simplification of example from Hartford, Lewis, Leyton-Brown, Taddy (2017): Deep IV: A Flexible 1/56 Approach for Counterfactual Prediction.

Unobserved variable X =desire for travel, affects both price (via airline algorithms) and seats sold.



Unobserved variable X =desire for travel, affects both price (via airline algorithms) and seats sold.



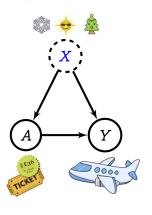
■ Desire for travel:

$$rac{oldsymbol{X}}{\mu} \sim \mathcal{N}(\mu, 0.01) \ \mu \sim \mathcal{U}\left\{-rac{1}{2}, 0, rac{1}{2}
ight\}$$

■ Price:

$$A = X + Z,$$
 $Z \sim \mathcal{N}(5, 0.04)$

Unobserved variable X =desire for travel, affects both price (via airline algorithms) and seats sold.



■ Desire for travel:

$$rac{oldsymbol{X}}{\mu} \sim \mathcal{N}(\mu, 0.01) \ \mu \sim \mathcal{U}\left\{-rac{1}{2}, 0, rac{1}{2}
ight\}$$

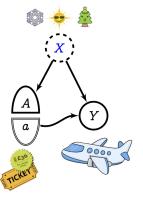
Price:

$$A = X + Z,$$
 $Z \sim \mathcal{N}(5, 0.04)$

■ Seats sold:

$$Y = 10 - A + 2X$$

Unobserved variable X =desire for travel, affects both price (via airline algorithms) and seats sold.



Desire for travel:

$$rac{oldsymbol{X}}{\mu} \sim \mathcal{N}(\mu, 0.01) \ \mu \sim \mathcal{U}\left\{-rac{1}{2}, 0, rac{1}{2}
ight\}$$

■ Price:

$$A = X + Z,$$
 $Z \sim \mathcal{N}(5, 0.04)$

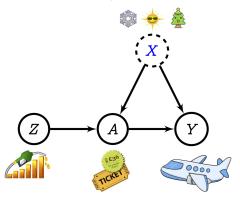
■ Seats sold:

$$Y=10-A+2X$$

Average treatment effect:

$$ext{ATE}(a) = \mathbb{E}[\,Y^{(a)}] = \int \left(10 - a + 2X
ight) dp(X) = 10 - a$$

Unobserved variable X =desire for travel, affects both price (via airline algorithms) and seats sold.



Desire for travel:

$$rac{oldsymbol{X}}{\mu \sim \mathcal{N}(\mu, 0.01)} \ \mu \sim \mathcal{U}\left\{-rac{1}{2}, 0, rac{1}{2}
ight\}$$

Price:

$$A = X + Z$$
, $Z \sim \mathcal{N}(5, 0.04)$

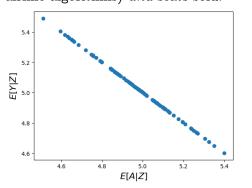
Seats sold:

$$Y=10-A+2X$$

Z is an instrument (cost of fuel). Condition on Z,

$$\mathbb{E}[Y|Z] = 10 - \mathbb{E}[A|Z] + 2\underbrace{\mathbb{E}[X|Z]}_{=0}$$

Unobserved variable X =desire for travel, affects both price (via airline algorithms) and seats sold.



■ Desire for travel:

$$egin{aligned} oldsymbol{X} &\sim \mathcal{N}(\mu, 0.01) \ \mu &\sim \mathcal{U}\left\{-rac{1}{2}, 0, rac{1}{2}
ight\} \end{aligned}$$

Price:

$$A = X + Z$$
, $Z \sim \mathcal{N}(5, 0.04)$

Seats sold:

$$Y=10-A+2X$$

Z is an instrument (cost of fuel). Condition on Z,

$$\mathbb{E}[Y|Z] = 10 - \mathbb{E}[A|Z] + 2\underbrace{\mathbb{E}[X|Z]}_{=0}$$

Instrumental variable regression

The Sveriges Riksbank Prize in Economic Sciences in Memory of Alfred Nobel 2021



© Nobel Prize Outreach. Photo: Paul Kennedy David Card Prize share: 1/2



© Nobel Prize Outreach. Photo: Risdon Photography Joshua D. Angrist Prize share: 1/4



© Nobel Prize Outreach. Photo: Paul Kennedy Guido W. Imbens Prize share: 1/4

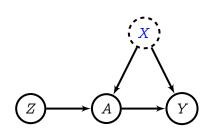
The Sveriges Riksbank Prize in Economic Sciences in Memory of Alfred Nobel 2021 was divided, one half awarded to David Card "for his empirical contributions to labour economics", the other half jointly to Joshua D. Angrist and Guido W. Imbens "for their methodological contributions to the analysis of causal relationships"

Instrumental variable regression with NN features

- X: unobserved confounder.
- A: treatment
- Y: outcome
- \blacksquare Z: instrument

Assumptions

$$egin{aligned} \mathbb{E}[X|Z] &= 0 \ Z \not\perp A \ (Y \perp\!\!\!\!\perp Z|A)_{G_{ar{A}}} \ Y &= \gamma^{ op} \phi_{ heta}(A) + X \end{aligned}$$

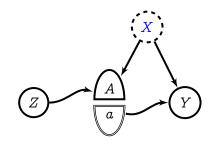


Instrumental variable regression with NN features

- X: unobserved confounder.
- A: treatment
- Y: outcome
- \blacksquare Z: instrument

Assumptions

$$egin{aligned} \mathbb{E}[X|Z] &= 0 \ Z \not\perp A \ (Y \perp\!\!\!\perp Z|A)_{G_{ar{A}}} \ Y &= \gamma^ op \phi_{ heta}(A) + X \end{aligned}$$



Average causal effect:

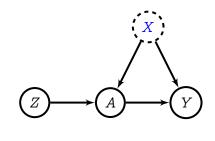
$$ext{ATE}(a) = \int \mathbb{E}(\left. Y \middle| oldsymbol{X}, a
ight) dp(oldsymbol{X}) = oldsymbol{\gamma}^ op oldsymbol{\phi}_{ heta}(a)$$

Instrumental variable regression with NN features

- *X*: unobserved confounder.
- A: treatment
- Y: outcome
- \blacksquare Z: instrument

Assumptions

$$egin{aligned} \mathbb{E}[X|Z] &= 0 \ Z \not\perp \!\!\! \perp A \ (Y \perp\!\!\! \perp Z|A)_{G_{ar{A}}} \ Y &= \gamma^{ op} \phi_{ heta}(A) + X \end{aligned}$$



IV regression: Condition both sides on Z,

$$\mathbb{E}[Y|Z] = \gamma^{ op} \mathbb{E}[\phi_{ heta}(A)|Z] + \underbrace{\mathbb{E}[X|Z]}_{=0}$$

Two-stage least squares for IV regression

Kernel features (NeurIPS 2019):



NN features (ICLR 2021):











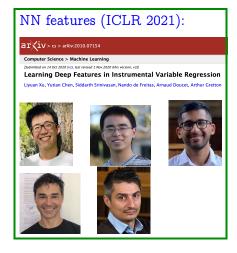


Code for NN and kernel IV methods: https://github.com/liyuan9988/DeepFeatureIV/

Two-stage least squares for IV regression

Kernel features (NeurIPS 2019):





Code for NN and kernel IV methods:

Stage 2 regression (IV): learn NN features $\phi_{\theta}(A)$ and linear layer γ to obtain Y with RR loss:

$$\mathbb{E}_{YZ}\left[(Y-\pmb{\gamma}^{ op}\mathbb{E}[\pmb{\phi}_{ heta}(A)|Z])^2
ight]+\lambda_2\|\pmb{\gamma}\|^2$$

Stage 2 regression (IV): learn NN features $\phi_{\theta}(A)$ and linear layer γ to obtain Y with RR loss:

$$\mathbb{E}_{YZ}\left[(Y-\pmb{\gamma}^{ op}\mathbb{E}[\pmb{\phi}_{ heta}(A)|Z])^2
ight]+\lambda_2\|\pmb{\gamma}\|^2$$

Stage 1 regression: learn NN features $\phi_{\zeta}(Z)$ and linear layer F:

$$\mathbb{E}[\phi_{ heta}(A)|Z] pprox F\phi_{\zeta}(Z)$$

with RR loss

$$\mathbb{E} \|\phi_{ heta}(A) - F\phi_{\zeta}(Z)\|^2 + \lambda_1 \|F\|_{HS}^2$$

Stage 2 regression (IV): learn NN features $\phi_{\theta}(A)$ and linear layer γ to obtain Y with RR loss:

$$\mathbb{E}_{YZ}\left[(Y-\pmb{\gamma}^{ op}\mathbb{E}[\pmb{\phi}_{ heta}(A)|Z])^2
ight]+\lambda_2\|\pmb{\gamma}\|^2$$

Stage 1 regression: learn NN features $\phi_{\zeta}(Z)$ and linear layer F:

$$\mathbb{E}[\phi_{ heta}(A)|Z]pprox F\phi_{\zeta}(Z)$$

with RR loss

$$\mathbb{E} \|\phi_{ heta}(A) - F\phi_{\zeta}(Z)\|^2 + \lambda_1 \|F\|_{HS}^2$$

Challenge: how to learn θ ?

Stage 2 regression (IV): learn NN features $\phi_{\theta}(A)$ and linear layer γ to obtain Y with RR loss:

$$\mathbb{E}_{YZ}\left[(Y-\pmb{\gamma}^{ op}\mathbb{E}[\pmb{\phi}_{ heta}(A)|Z])^2
ight]+\lambda_2\|\pmb{\gamma}\|^2$$

Stage 1 regression: learn NN features $\phi_{\zeta}(Z)$ and linear layer F:

$$\mathbb{E}[\phi_{ heta}(A)|Z]pprox F\phi_{\zeta}(Z)$$

with RR loss

$$\mathbb{E} \|\phi_{ heta}(A) - F\phi_{\zeta}(Z)\|^2 + \lambda_1 \|F\|_{HS}^2$$

Challenge: how to learn θ ?

From Stage 2 regression?

Stage 2 regression (IV): learn NN features $\phi_{\theta}(A)$ and linear layer γ to obtain Y with RR loss:

$$\mathbb{E}_{YZ}\left[(Y-\pmb{\gamma}^{ op}\mathbb{E}[\pmb{\phi}_{ heta}(A)|Z])^2
ight]+\lambda_2\|\pmb{\gamma}\|^2$$

Stage 1 regression: learn NN features $\phi_{\zeta}(Z)$ and linear layer F:

$$\mathbb{E}[\phi_{ heta}(A)|Z]pprox F\phi_{\zeta}(Z)$$

with RR loss

$$\mathbb{E} \|\phi_{ heta}(A) - F\phi_{\zeta}(Z)\|^2 + \lambda_1 \|F\|_{HS}^2$$

Challenge: how to learn θ ?

From Stage 2 regression?

...which requires $\mathbb{E}[\phi_{\theta}(A)|Z]$ from Stage 1 regression

Stage 2 regression (IV): learn NN features $\phi_{\theta}(A)$ and linear layer γ to obtain Y with RR loss:

$$\mathbb{E}_{YZ}\left[(Y-\pmb{\gamma}^{ op}\mathbb{E}[\pmb{\phi}_{ heta}(A)|Z])^2
ight]+\lambda_2\|\pmb{\gamma}\|^2$$

Stage 1 regression: learn NN features $\phi_{\zeta}(Z)$ and linear layer F:

$$\mathbb{E}[\phi_{ heta}(A)|Z]pprox F\phi_{\zeta}(Z)$$

with RR loss

$$\mathbb{E} \|\phi_{ heta}(A) - F\phi_{\zeta}(Z)\|^2 + \lambda_1 \|F\|_{HS}^2$$

Challenge: how to learn θ ?

From Stage 2 regression?

...which requires $\mathbb{E}[\phi_{\theta}(A)|Z]$ from Stage 1 regression

...which requires $\phi_{\theta}(A)$... which requires θ ...

Stage 2 regression (IV): learn NN features $\phi_{\theta}(A)$ and linear layer γ to obtain Y with RR loss:

$$\mathbb{E}_{YZ}\left[(Y-\pmb{\gamma}^{ op}\mathbb{E}[\pmb{\phi}_{ heta}(A)|Z])^2
ight]+\lambda_2\|\pmb{\gamma}\|^2$$

Stage 1 regression: learn NN features $\phi_{\zeta}(Z)$ and linear layer F:

$$\mathbb{E}[\phi_{ heta}(A)|Z] pprox F\phi_{\zeta}(Z)$$

with RR loss

$$\mathbb{E} \|\phi_{ heta}(A) - F\phi_{\zeta}(Z)\|^2 + \lambda_1 \|F\|_{HS}^2$$

Challenge: how to learn θ ?

From Stage 2 regression?

...which requires $\mathbb{E}[\phi_{\theta}(A)|Z]$ from Stage 1 regression

...which requires $\phi_{\theta}(A)$... which requires θ ...

Use the linear final layers! (i.e. γ and F)

Xu, Chen, Srinivasan, De Freitas, Doucet, G. (2021) Learning Deep Features in Instrumental Variable Regresion 27/56

Stage 1 regression: learn NN features $\phi_{\zeta}(Z)$ and linear layer F:

$$\mathbb{E}[\phi_{ heta}(A)|Z] pprox F\phi_{\zeta}(Z)$$

with RR loss

$$\mathbb{E}\left[\|\phi_{ heta}(A) - F\phi_{\zeta}(Z)\|^2
ight] + \lambda_1 \|F\|_{HS}^2$$

Stage 1 regression: learn NN features $\phi_{\zeta}(Z)$ and linear layer F:

$$\mathbb{E}[\phi_{ heta}(A)|Z] pprox F\phi_{\zeta}(Z)$$

with RR loss

$$\mathbb{E}\left[\|\phi_{\theta}(A) - {\color{red}F}\phi_{\zeta}(Z)\|^2\right] + \lambda_1\|{\color{red}F}\|_{HS}^2$$

 $\hat{F}_{\theta,\zeta}$ in closed form wrt $\phi_{\theta},\phi_{\zeta}$:

$$egin{aligned} \hat{F}_{ heta,\zeta} &= C_{AZ}(C_{ZZ} + \lambda_1 I)^{-1} & C_{AZ} &= \mathbb{E}[\phi_{ heta}(A)\phi_{\zeta}^{ op}(Z)] \ C_{ZZ} &= \mathbb{E}[\phi_{\zeta}(Z)\phi_{\zeta}^{ op}(Z)] \end{aligned}$$

Stage 1 regression: learn NN features $\phi_{\zeta}(Z)$ and linear layer F:

$$\mathbb{E}[\phi_{ heta}(A)|Z] pprox F\phi_{\zeta}(Z)$$

with RR loss

$$\mathbb{E}\left[\|\phi_{ heta}(A) - F\phi_{\zeta}(Z)\|^2
ight] + \lambda_1 \|F\|_{HS}^2$$

 $\hat{F}_{\theta,\zeta}$ in closed form wrt $\phi_{\theta},\phi_{\zeta}$:

$$egin{aligned} \hat{F}_{ heta,\zeta} &= C_{AZ}(C_{ZZ} + \lambda_1 I)^{-1} & C_{AZ} &= \mathbb{E}[\phi_{ heta}(A)\phi_{\zeta}^{ op}(Z)] \ C_{ZZ} &= \mathbb{E}[\phi_{\zeta}(Z)\phi_{\zeta}^{ op}(Z)] \end{aligned}$$

Plug $\hat{F}_{\theta,\zeta}$ into S1 loss, take gradient steps for ζ (...but not θ ...)

Stage 2 regression (IV): learn NN features $\phi_{\theta}(A)$ and linear layer γ to obtain Y with RR loss:

$$\mathcal{L}_2(\gamma, heta) = \mathbb{E}_{YZ}\left[(Y - \gamma^ op \mathbb{E}[\phi_{ heta}(A)|Z])^2
ight] + \lambda_2 \|\gamma\|^2$$

Stage 2 regression (IV): learn NN features $\phi_{\theta}(A)$ and linear layer γ to obtain Y with RR loss:

$$egin{aligned} \mathcal{L}_2(\gamma, heta) &= \mathbb{E}_{YZ}\left[(Y - \gamma^ op \mathbb{E}[\phi_{ heta}(A)|Z])^2
ight] + \lambda_2 \|\gamma\|^2 \ &= \mathbb{E}_{YZ}[(Y - \gamma^ op rac{\hat{F}_{ heta,\zeta}\phi_{\zeta}(Z)}{ ext{Stare 1}})^2] + \lambda_2 \|\gamma\|^2 \end{aligned}$$

Stage 2 regression (IV): learn NN features $\phi_{\theta}(A)$ and linear layer γ to obtain Y with RR loss:

$$egin{aligned} \mathcal{L}_2(\gamma, heta) &= \mathbb{E}_{YZ}\left[(\,Y - \gamma^ op \mathbb{E}[\phi_ heta(A)|Z])^2
ight] + \lambda_2 \|\gamma\|^2 \ &= \mathbb{E}_{\,YZ}[(\,Y - \gamma^ op \hat{F}_{\, heta,\zeta}\phi_\zeta(Z))^2] + \lambda_2 \|\gamma\|^2 \end{aligned}$$

 $\hat{\gamma}_{\theta}$ in closed form wrt ϕ_{θ} :

$$egin{aligned} \hat{\gamma}_{ heta} &:= \widetilde{C}_{YA|Z} (\widetilde{C}_{AA|Z} + \lambda_2 I)^{-1} \qquad \widetilde{C}_{YA|Z} = \mathbb{E} \left[Y \ [\hat{F}_{ heta,\zeta} arphi_{\zeta}(Z)]^{ op}
ight] \ \widetilde{C}_{AA|Z} &= \mathbb{E} \left[[\hat{F}_{ heta,\zeta} arphi_{\zeta}(Z)] \ [\hat{F}_{ heta,\zeta} arphi_{\zeta}(Z)]^{ op}
ight] \end{aligned}$$

Stage 2 regression (IV): learn NN features $\phi_{\theta}(A)$ and linear layer γ to obtain Y with RR loss:

$$egin{aligned} \mathcal{L}_2(\gamma, heta) &= \mathbb{E}_{YZ}\left[(\,Y - \gamma^ op \mathbb{E}[\phi_ heta(A)|Z])^2
ight] + \lambda_2 \|\gamma\|^2 \ &= \mathbb{E}_{\,YZ}[(\,Y - \gamma^ op \hat{F}_{\, heta,\zeta}\phi_\zeta(Z))^2] + \lambda_2 \|\gamma\|^2 \end{aligned}$$

 $\hat{\gamma}_{\theta}$ in closed form wrt ϕ_{θ} :

$$egin{aligned} \hat{\gamma}_{ heta} &:= \widetilde{C}_{YA|Z} (\widetilde{C}_{AA|Z} + \lambda_2 I)^{-1} \qquad \widetilde{C}_{YA|Z} = \mathbb{E} \left[Y \ [\hat{m{F}}_{ heta,\zeta} m{arphi}_{\zeta}(Z)]^{ op}
ight] \ \widetilde{C}_{AA|Z} &= \mathbb{E} \left[[\hat{m{F}}_{ heta,\zeta} m{arphi}_{\zeta}(Z)] \ [\hat{m{F}}_{ heta,\zeta} m{arphi}_{\zeta}(Z)]^{ op}
ight] \end{aligned}$$

From linear final layers in Stages 1,2:

Learn $\phi_{\theta}(A)$ by plugging $\hat{\gamma}_{\theta}$ into S2 loss, taking gradient steps for θ

Stage 2 regression (IV): learn NN features $\phi_{\theta}(A)$ and linear layer γ to obtain Y with RR loss:

$$egin{aligned} \mathcal{L}_2(\gamma, heta) &= \mathbb{E}_{YZ}\left[(\,Y - \gamma^ op \mathbb{E}[\phi_ heta(A)|Z])^2
ight] + \lambda_2 \|\gamma\|^2 \ &= \mathbb{E}_{\,YZ}[(\,Y - \gamma^ op \hat{m{F}}_{ heta,\zeta}\phi_\zeta(Z))^2] + \lambda_2 \|\gamma\|^2 \end{aligned}$$

 $\hat{\gamma}_{\theta}$ in closed form wrt ϕ_{θ} :

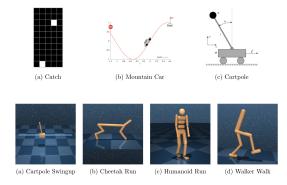
$$egin{aligned} \hat{\gamma}_{ heta} &:= \widetilde{C}_{YA|Z} (\widetilde{C}_{AA|Z} + \lambda_2 I)^{-1} \qquad \widetilde{C}_{YA|Z} = \mathbb{E} \left[Y \left[\hat{F}_{ heta,\zeta} arphi_{\zeta} (Z)
ight]^{ op}
ight] \ \widetilde{C}_{AA|Z} &= \mathbb{E} \left[\left[\hat{F}_{ heta,\zeta} arphi_{\zeta} (Z)
ight] \left[\hat{F}_{ heta,\zeta} arphi_{\zeta} (Z)
ight]^{ op}
ight] \end{aligned}$$

From linear final layers in Stages 1,2:

Learn $\phi_{\theta}(A)$ by plugging $\hat{\gamma}_{\theta}$ into S2 loss, taking gradient steps for θ but $\hat{\zeta}$ changes with θ

...so alternate first and second stages until convergence.

Neural IV in reinforcement learning



Policy evaluation: want Q-value:

$$Q^{\pi}(s,a) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t R_t \middle| S_0 = s, A_0 = a
ight]$$

for policy $\pi(A|S=s)$.

Osband et al (2019). Behaviour suite for reinforcement learning.https://github.com/deepmind/bsuite

Tassa et al. (2020). dm_control:Software and tasks for continuous control.

30/56
https://github.com/deepmind/dm_control

Application of IV: reinforcement learning

Q value is a minimizer of Bellman loss

$$\mathcal{L}_{ ext{Bellman}} = \mathbb{E}_{\mathit{SAR}}\left[\left(R + \gamma[\mathbb{E}\left[\left.Q^{\pi}(S',A')\middle|S,A
ight] - \left.Q^{\pi}(S,A)
ight)^2
ight].$$

Corresponds to "IV-like" problem

$$\mathcal{L}_{ ext{Bellman}} = \mathbb{E}_{\,YZ}\left[\left(\,Y - \mathbb{E}[f(X)|Z]
ight)^2
ight]$$

with

$$egin{aligned} Y &= R, \ X &= (S', A', S, A) \ Z &= (S, A), \ f_0(X) &= Q^\pi(s,a) - \gamma Q^\pi(s',a') \end{aligned}$$

RL experiments and data:

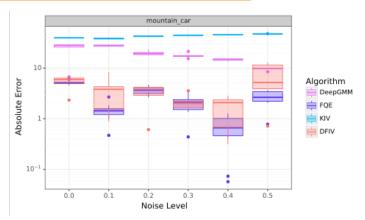
https://github.com/liyuan9988/IVOPEwithACME

Bradtke and Barto (1996). Linear least-squares algorithms for temporal difference learning.

Xu, Chen, Srinivasan, De Freitas, Doucet, G. (2021)

Chen, Xu, Gulcehre, Le Paine, G, De Freitas, Doucet (2022). On Instrumental Variable Regression 1956 Deep Offline Policy Evaluation.

Results on mountain car problem



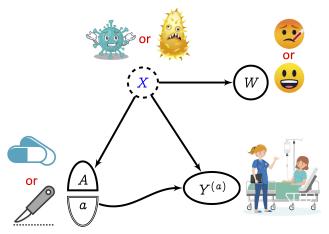
Good performance compared with FQE.

Warning: IV assumption can fail when regression underfits. See papers for details.

Xu, Chen, Srinivasan, De Freitas, Doucet, G. (2021)
Chen, Xu, Gulcehre, Le Paine, G, De Freitas, Doucet (2022). On Instrumental Variable Regression 2956
Deep Offline Policy Evaluation.

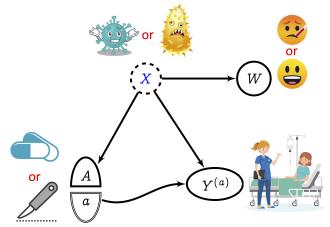
...but seriously, what if there are hidden confounders?

We record symptom W, not disease X



- P(W = fever|X = mild) = 0.2
- P(W = fever|X = severe) = 0.8

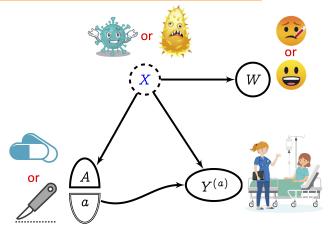
We record symptom W, not disease X



- P(W = fever|X = mild) = 0.2
- P(W = fever|X = severe) = 0.8

Could we just write: $P(Y^{(a)}) \stackrel{?}{=} \sum_{w \in \{0,1\}} \mathbb{E}[Y|a,w] p(w)$

We record symptom W, not disease X



Wrong recommendation made:

- \blacksquare $\sum_{w \in \{0,1\}} \mathbb{E}[\text{cured}|\text{surgery}, w] p(w) = 0.73 \quad (\neq 0.75)$

Correct answer impossible without observing X

Proxy causal learning (negative controls)

Causal effect estimation, with hidden covariates X:

■ Use proxy variables (negative controls)

Applications: effect of actions under

- privacy constraints (email, ads, DMA)
- data gathering constraints (edge computing)
- fundamental limitations (preferences, state of mind)

Proxy causal learning (negative controls)

Causal effect estimation, with hidden covariates X:

■ Use proxy variables (negative controls)

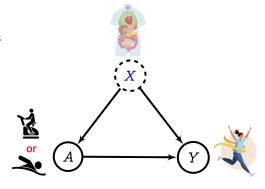
Applications: effect of actions under

- privacy constraints (email, ads, DMA)
- data gathering constraints (edge computing)
- fundamental limitations (preferences, state of mind)

Don't meet your heroes model your hidden variables!

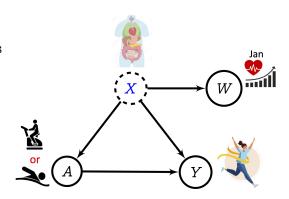
Unobserved X with (possibly) complex nonlinear effects on A, Y

- X: true physical status
- A: exercise regimes
- Y: fitness goal



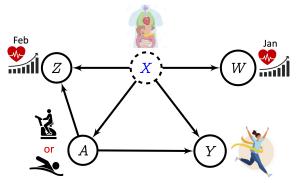
Unobserved X with (possibly) complex nonlinear effects on A, Y

- X: true physical status
- A: exercise regimes
- Y: fitness goal
- W: health readings before A



Unobserved X with (possibly) complex nonlinear effects on A, Y

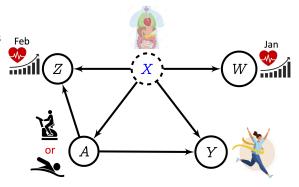
- \blacksquare X: true physical status
- A: exercise regimes
- \blacksquare Y: fitness goal
- W: health readings before A
- Z: health readings after A



Unobserved X with (possibly) complex nonlinear effects on A, Y

In this example:

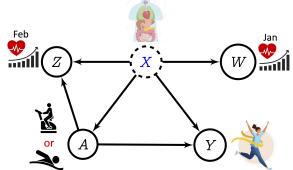
- \blacksquare X: true physical status
- A: exercise regimes
- Y: fitness goal
- W: health readings before A
- Z: health readings after A



 \implies Can recover $\mathbb{E}(Y^{(a)})$ from observational data

Unobserved X with (possibly) complex nonlinear effects on A, Y

- \blacksquare X: true physical status
- A: exercise regimes
- Y: fitness goal
- W: health readings before A
- Z: health readings after A



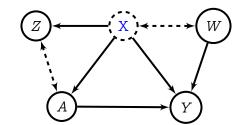
- \implies Can recover $\mathbb{E}(Y^{(a)})$ from observational data
- ⇒ More usefully: evaluate novel, on-device policy:

$$\mathbb{E}(Y^{(\pi(A|X))})$$

Proxy variables: general setting

Unobserved X with (possibly) complex nonlinear effects on A, Y The definitions are:

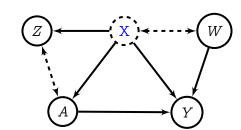
- X: unobserved confounder.
- *A*: treatment
- *Y*: outcome
- \blacksquare Z: treatment proxy
- W outcome proxy



Proxy variables: general setting

Unobserved X with (possibly) complex nonlinear effects on A, Y. The definitions are:

- *X*: unobserved confounder.
- A: treatment
- *Y*: outcome
- \blacksquare Z: treatment proxy
- W outcome proxy



Structural assumptions:

$$W \perp \!\!\!\perp (Z, A)|X$$

 $Y \perp \!\!\!\perp Z|(A, X)$

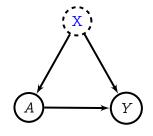
Miao, Geng, Tchetgen Tchetgen (2018): Identifying causal effects with proxy variables of an unmeasured confounder.

37/56

Why proxy variables? A simple proof

The definitions are:

- X: unobserved confounder.
- A: treatment
- *Y*: outcome



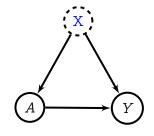
If X were observed,

$$\underbrace{P(Y^{(a)})}_{d_{ai} imes 1} := \sum_{i=1}^{d_{ai}} P(Y|\mathbf{x}_i, a) P(\mathbf{x}_i)$$

Why proxy variables? A simple proof

The definitions are:

- X: unobserved confounder.
- A: treatment
- *Y*: outcome



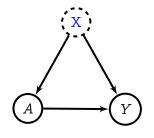
If X were observed,

$$\underbrace{P(Y^{(a)})}_{d_y \times 1} := \sum_{i=1}^{d_x} P(Y|x_i, a) P(x_i) = \underbrace{P(Y|X, a) P(X)}_{d_y \times d_x} \underbrace{P(X|X, a) P(X)}_{d_x \times 1}$$

Why proxy variables? A simple proof

The definitions are:

- X: unobserved confounder.
- A: treatment
- *Y*: outcome



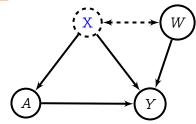
If X were observed,

$$\underbrace{P(Y^{(a)})}_{d_y imes 1} := \sum_{i=1}^{d_x} P(Y|x_i, a) P(x_i) = \underbrace{P(Y|X, a)}_{d_y imes d_x} \underbrace{P(X)}_{d_x imes 1}$$

Goal: "get rid of the blue" X

The definitions are:

- X: unobserved confounder.
- A: treatment
- *Y*: outcome
- W: outcome proxy

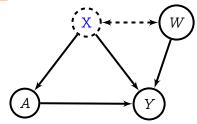


For each a, if we could solve:

$$\underbrace{P(Y|X,a)}_{d_y imes d_x} = \underbrace{H_{w,a}}_{d_y imes d_w} \underbrace{P(W|X)}_{d_w imes d_x}$$

The definitions are:

- X: unobserved confounder.
- A: treatment
- *Y*: outcome
- W: outcome proxy



For each a, if we could solve:

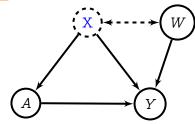
$$\underbrace{P(Y|X,a)}_{d_y imes d_x} = \underbrace{H_{w,a}}_{d_y imes d_w} \underbrace{P(W|X)}_{d_w imes d_x}$$

.....then

$$P(Y^{(a)}) = P(Y|X,a)P(X)$$

The definitions are:

- X: unobserved confounder.
- A: treatment
- *Y*: outcome
- W: outcome proxy



For each a, if we could solve:

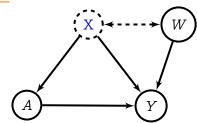
$$\underbrace{P(Y|X,a)}_{d_y imes d_x} = \underbrace{H_{w,a}}_{d_y imes d_w} \underbrace{P(W|X)}_{d_w imes d_x}$$

.....then

$$P(Y^{(a)}) = P(Y|X, a)P(X)$$
$$= H_{w,a}P(W|X)P(X)$$

The definitions are:

- *X*: unobserved confounder.
- A: treatment
- *Y*: outcome
- W: outcome proxy



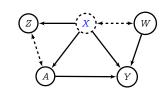
For each a, if we could solve:

$$\underbrace{P(\mathit{Y}|\mathit{X},\mathit{a})}_{\mathit{d_{\mathit{y}}} imes \mathit{d_{\mathit{x}}}} = \underbrace{\mathit{H_{w,a}}}_{\mathit{d_{\mathit{y}}} imes \mathit{d_{w}}} \underbrace{P(\mathit{W}|\mathit{X})}_{\mathit{d_{\mathit{w}}} imes \mathit{d_{x}}}$$

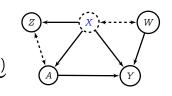
.....then

$$egin{aligned} P(Y^{(a)}) &= P(Y|X,a)P(X) \ &= H_{w,a}P(W|X)P(X) \ &= H_{w,a}P(W) \end{aligned}$$

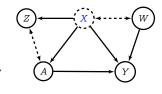
$$P(Y|X,a) = H_{w,a}P(W|X)$$



$$P(Y|X,a) \underbrace{p(X|Z,a)}_{d_x imes d_z} = H_{w,a} P(W|X) \underbrace{p(X|Z,a)}_{d_x imes d_z}$$



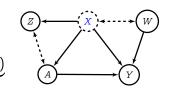
$$P(Y|X,a) \underbrace{p(X|Z,a)}_{d_x imes d_z} = H_{w,a} P(W|X) \underbrace{p(X|Z,a)}_{d_x imes d_z}$$



Because
$$W \perp \!\!\!\perp (Z,A)|X$$
,

$$P(W|X)p(X|Z,a) = P(W|Z,a)$$

$$P(Y|X,a)\underbrace{p(X|Z,a)}_{d_x \times d_z} = H_{w,a}P(W|X)\underbrace{p(X|Z,a)}_{d_x \times d_z}$$



Because
$$W \perp \!\!\!\perp (Z, A)|X$$
,

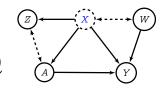
$$P(W|X)p(X|Z,a) = P(W|Z,a)$$

Because
$$Y \perp \!\!\!\perp Z | (A, X)$$
,

$$P(Y|X,a)p(X|Z,a) = P(Y|Z,a)$$

From last slide,

$$P(Y|X,a) \underbrace{p(X|Z,a)}_{d_x \times d_z} = H_{w,a} P(W|X) \underbrace{p(X|Z,a)}_{d_x \times d_z}$$



Because
$$W \perp \!\!\!\perp (Z, A)|X$$
,

$$P(W|X)p(X|Z,a) = P(W|Z,a)$$

Because
$$Y \perp \!\!\!\perp Z | (A, X)$$
,

$$P(Y|X,a)p(X|Z,a) = P(Y|Z,a)$$

Solve for $H_{w,a}$:

$$P(Y|Z,a) = H_{w,a}P(W|Z,a)$$

Everything observed!

Proxy/Negative Control Methods in the Real World

Unobserved confounders: proxy methods

Kernel features (ICML 2021):













NN features (NeurIPS 2021):







Code for NN and kernel proxy methods:

https://github.com/liyuan9988/DeepFeatureProxyVariable/

Unobserved confounders: proxy methods

Kernel features (ICML 2021):





Code for NN and kernel proxy methods:

https://github.com/liyuan9988/DeepFeatureProxyVariable/

Road map: NN proxy learning

We'll proceed as follows:

- Proxy relation for continuous variables
- Loss function for deep proxy learning
- Define primary (ridge) regression with this loss
- Define secondary (ridge) regression as input to primary

If X were observed, we would write (dose-response curve)

$$\mathbb{E}(Y^{(a)}) = \int_x \mathbb{E}(Y|a,x)p(x)dx.$$

....but we do not observe X.

If X were observed, we would write (dose-response curve)

$$\mathbb{E}(Y^{(a)}) = \int_x \mathbb{E}(Y|a,x)p(x)dx.$$

....but we do not observe X.

Main theorem: Assume we solved for link function:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

- "Primary" $\mathbb{E}(Y|a,z)$, "secondary" $\mathbb{E}_{W|a,z}$ linked by h_y
- All variables observed, X not seen or modeled.

Fredholm equation of first kind. Link existence requires \diamondsuit , identification of ATE requires \triangle (and further technical assumptions) [XKG: Asspumption 2, Prop. 1,Corr. 1; Deaner]

$$\mathbb{E}[f(X)|A=a,Z=z]=0,\ \forall (z,a)\iff f(X)=0,\ \mathbb{P}_X \ \mathrm{a.s.} \ \triangle$$
 $\mathbb{E}[f(X)|A=a,W=w]=0,\ \forall (w,a)\iff f(X)=0,\ \mathbb{P}_X \ \mathrm{a.s.} \ \diamondsuit$

If X were observed, we would write (dose-response curve)

$$\mathbb{E}(Y^{(a)}) = \int_x \mathbb{E}(Y|a,x)p(x)dx.$$

....but we do not observe X.

Main theorem: Assume we solved for link function:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z} h_y(W,a)$$

- "Primary" $\mathbb{E}(Y|a,z)$, "secondary" $\mathbb{E}_{W|a,z}$ linked by h_y
- \blacksquare All variables observed, X not seen or modeled.

Dose-response curve via p(w):

$$\mathbb{E}(\,Y^{(a)}) = \int_w h_y(a,w) p(w) dw$$

If X were observed, we would write (dose-response curve)

$$\mathbb{E}(Y^{(a)}) = \int_x \mathbb{E}(Y|a,x)p(x)dx.$$

....but we do not observe X.

Main theorem: Assume we solved for link function:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

- "Primary" $\mathbb{E}(Y|a,z)$, "secondary" $\mathbb{E}_{W|a,z}$ linked by h_y
- \blacksquare All variables observed, X not seen or modeled.

Dose-response curve via p(w):

$$\mathbb{E}(\mathit{Y}^{(a)}) = \int_{w} \mathit{h}_{y}(a,w) \mathit{p}(w) \mathit{d}w$$

Challenge: need a loss function for h_v

Primary loss function for $h_y(w, a)$

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Primary loss function:

$$\hat{h}_{y} = rg \min_{h_{y}} \mathbb{E}_{Y,A,Z} \left(Y - \mathbb{E}_{W|A,Z} h_{y}(W,A) \right)^{2}$$

Why?

Primary loss function for $h_y(w, a)$

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Primary loss function:

$$\hat{h}_{y} = rg\min_{h_{y}} \mathbb{E}_{Y,A,Z} \left(Y - \mathbb{E}_{W|A,Z} h_{y}(W,A) \right)^{2}$$

Why?

$$f^*(a,z) = \mathbb{E}(\,Y|\,a,z) ext{ solves}
onumber \ \, rgmin_f \mathbb{E}_{\,Y,A,Z} \, (\,Y-f(A,Z))^2$$

Primary loss function for $h_y(w, a)$

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Primary loss function:

$$\hat{h}_{y} = rg \min_{h_{y}} \mathbb{E}_{Y,A,Z} \left(Y - \mathbb{E}_{W|A,Z} h_{y}(W,A) \right)^{2}$$

Why?

$$f^*(a,z) = \mathbb{E}(\,Y|\,a,z) ext{ solves}
onumber rgmin } \mathbb{E}_{\,Y,A,Z} \, (\,Y - f(A,Z))^2$$

...and by the proxy model above,

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z} h_y(W,a)$$

Deaner (2021). Mastouri, Zhu, Gultchin, Korba, Silva, Kusner, G., Muandet (2021). Xu, Kanagawa, G. (2021).

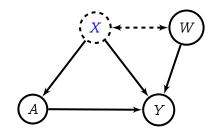
NN for link $h_y(a, w)$

The link function is a function of two arguments

$$h_y(a,w) = \gamma^ op \left[arphi_{ heta}(w) \otimes arphi_{\xi}(a)
ight] = \gamma^ op \left[egin{array}{c} arphi_{ heta,1}(w) arphi_{\xi,1}(a) \ arphi_{ heta,1}(w) arphi_{\xi,2}(a) \ dots \ dots \ arphi_{ heta,2}(w) arphi_{\xi,1}(a) \ dots \ dots \ \end{array}
ight]$$

Assume we have:

- lacksquare output proxy NN features $arphi_{ heta}(w)$
- lacksquare treatment NN features $arphi_{\xi}(a)$
- linear final layer γ
 (argument of feature map indicates feature space)



NN for link $h_y(a, w)$

The link function is a function of two arguments

$$h_y(a,w) = \gamma^ op \left[arphi_ heta(w) \otimes arphi_\xi(a)
ight]$$

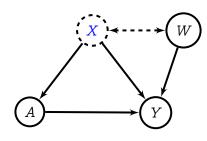
Assume we have:

- lacksquare output proxy NN features $\varphi_{\theta}(w)$
- lacksquare treatment NN features $arphi_{\xi}(a)$
- In linear final layer γ (argument of feature map indicates feature space)

Questions:

- Why feature map $\varphi_{\theta}(w) \otimes \varphi_{\xi}(a)$?
- Why final linear layer γ ?

Both are necessary (next slide)!



NN ridge regression for $h_y(w, a)$

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Primary regression:

$$\hat{h}_y = rg\min_{h_y} \mathbb{E}_{Y,A,Z} \left(Y - \mathbb{E}_{oldsymbol{W}|A,Z} h_y(oldsymbol{W},A)
ight)^2 + \lambda_2 \|\gamma\|^2$$

NN ridge regression for $h_y(w, a)$

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Primary regression:

$$\hat{h}_{y} = \arg\min_{h_{y}} \mathbb{E}_{Y,A,Z} \left(\left. Y - \mathbb{E}_{\boldsymbol{W}|\boldsymbol{A},\boldsymbol{Z}} h_{y}(\boldsymbol{W},\boldsymbol{A}) \right)^{2} + \lambda_{2} \| \gamma \|^{2} \right.$$

How to get conditional expectation $\mathbb{E}_{W|a,z}h_y(W,a)$?

Density estimation for p(W|a, z)? Sample from p(W|a, z)?

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Primary regression:

$$\hat{h}_y = \arg\min_{h_y} \mathbb{E}_{Y,A,Z} \left(\left. Y - \mathbb{E}_{\left. \boldsymbol{W} \right| A,Z} h_y (\left. \boldsymbol{W},A \right) \right)^2 + \lambda_2 \| \gamma \|^2$$

Recall link function

$$h_y(extit{W}, a) = \left[\gamma^ op (arphi_ heta(extit{W}) \otimes arphi_\xi(a))
ight]$$

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Primary regression:

$$\hat{h}_y = \arg\min_{h_y} \mathbb{E}_{Y,A,Z} \left(\left. Y - \mathbb{E}_{\left. W \right| A,Z} h_y (\left. W \right,A) \right)^2 + \lambda_2 \|\gamma\|^2$$

Recall link function

Xu, Kanagawa, G. (2021).

$$\mathbb{E}_{W|a,z} \; h_y(\hspace{.05cm} W,\hspace{.05cm} a) = \hspace{.05cm} \mathbb{E}_{\hspace{.05cm} W|a,z} \; \left[\gamma^{ op} (\hspace{.05cm} arphi_{ heta}(\hspace{.05cm} W) \otimes arphi_{\xi}(\hspace{.05cm} a))
ight]$$

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Primary regression:

$$\hat{h}_y = \arg\min_{h_y} \mathbb{E}_{Y,A,Z} \left(\left. Y - \mathbb{E}_{\left. W \right| A,Z} h_y (\left. W \right, A) \right)^2 + \lambda_2 \| \gamma \|^2$$

Recall link function

$$egin{aligned} \mathbb{E}_{W|a,z} \; h_y(\,W,\,a) &= \; \mathbb{E}_{W|a,z} \; \left[\gamma^ op \left(arphi_ heta(\,W) \otimes arphi_\xi(a)
ight)
ight] \ &= \gamma^ op \left(\mathbb{E}_{W|a,z} \left[arphi_ heta(\,W)
ight] \otimes arphi_\xi(a)
ight) \ & ext{cond. feat. mean} \end{aligned}$$

(this is why linear γ and feature map $\varphi_{\theta}(w) \otimes \varphi_{\xi}(a)$)

Goal:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z}h_y(W,a)$$

Primary regression:

$$\hat{h}_y = \arg\min_{h_y} \mathbb{E}_{Y,A,Z} \left(\left. Y - \mathbb{E}_{\left. W \right| A,Z} h_y (\left. W \right, A) \right)^2 + \lambda_2 \| \gamma \|^2$$

Recall link function

$$egin{aligned} \mathbb{E}_{W|a,z} \; h_y(\,W,\,a) &= \; \mathbb{E}_{W|a,z} \; \left[\gamma^ op \left(arphi_ heta(\,W) \otimes arphi_\xi(a)
ight)
ight] \ &= \gamma^ op \left(\mathbb{E}_{W|a,z} \left[arphi_ heta(\,W)
ight] \otimes arphi_\xi(a)
ight) \ & ext{cond. feat. mean} \end{aligned}$$

Ridge regression (again!)

$$\mathbb{E}_{W|a,z}arphi_{ heta}(W)=\hat{F}_{ heta,\zeta}arphi_{\zeta}(a,z)$$

NN ridge regression for $\mathbb{E}_{W|a,z}\varphi_{\theta}(W)$

Secondary regression: learn NN features $\varphi_{\zeta}(Z)$ and linear layer F:

$$\mathbb{E}_{W|a,z}arphi_{ heta}(W)=\hat{F}_{ heta,\zeta}arphi_{\zeta}(a,z)$$

with RR loss

$$\mathbb{E}_{W,A,Z} \left\| arphi_{ heta}(W) - rac{oldsymbol{F}}{oldsymbol{F}} arphi_{\zeta}(A,Z)
ight\|^2 + \lambda_1 \|rac{oldsymbol{F}}{oldsymbol{F}} \|^2$$

 $\hat{F}_{\theta,\zeta}$ in closed form wrt $\varphi_{\theta}, \varphi_{\zeta}$.

NN ridge regression for $\mathbb{E}_{W|a,z}\varphi_{\theta}(W)$

Secondary regression: learn NN features $\varphi_{\zeta}(Z)$ and linear layer F:

$$\mathbb{E}_{W|a,z}arphi_{ heta}({W})=\hat{F}_{ heta,\zeta}arphi_{\zeta}(a,z)$$

with RR loss

$$\mathbb{E}_{W,A,Z} \left\| arphi_{ heta}(W) - rac{oldsymbol{F}}{oldsymbol{F}} arphi_{\zeta}(A,Z)
ight\|^2 + \lambda_1 \|rac{oldsymbol{F}}{oldsymbol{F}} \|^2$$

 $\hat{F}_{\theta,\zeta}$ in closed form wrt $\varphi_{\theta}, \varphi_{\zeta}$.

Plug $\hat{F}_{\theta,\zeta}$ into S1 loss, backprop through Cholesky for ζ (...not θ ...why not?)

```
Solve for \theta, \xi, \zeta:
```

Repeat until convergence:

■ Secondary: Solve for $\hat{F}_{\theta,\zeta}$, then gradient steps on ζ (backprop through Cholesky)

```
Solve for \theta, \xi, \zeta:
```

Repeat until convergence:

- Secondary: Solve for $\hat{F}_{\theta,\zeta}$, then gradient steps on ζ (backprop through Cholesky)
- Primary: Solve for $\hat{\gamma}$ in terms of $\hat{F}_{\theta,\zeta}\varphi_{\zeta}(A,Z)$ and $\varphi_{\xi}(A)$

```
Solve for \theta, \xi, \zeta:
```

Repeat until convergence:

- Secondary: Solve for $\hat{F}_{\theta,\zeta}$, then gradient steps on ζ (backprop through Cholesky)
- Primary: Solve for $\hat{\gamma}$ in terms of $\hat{F}_{\theta,\zeta}\varphi_{\zeta}(A,Z)$ and $\varphi_{\xi}(A)$
- Primary: Gradient steps on θ , ξ (backprop through Cholesky)
 - $\hat{F}_{\theta,\zeta}$ remains optimal wrt current φ_{θ} .

```
Solve for \theta, \xi, \zeta:
```

Repeat until convergence:

- Secondary: Solve for $\hat{F}_{\theta,\zeta}$, then gradient steps on ζ (backprop through Cholesky)
- Primary: Solve for $\hat{\gamma}$ in terms of $\hat{F}_{\theta,\zeta}\varphi_{\zeta}(A,Z)$ and $\varphi_{\xi}(A)$
- Primary: Gradient steps on θ, ξ (backprop through Cholesky)
 - $\hat{F}_{\theta,\zeta}$ remains optimal wrt current φ_{θ} .

Iterate between updates of θ , ξ and ζ

Solve for θ, ξ, ζ :

Repeat until convergence:

- Secondary: Solve for $\hat{F}_{\theta,\zeta}$, then gradient steps on ζ (backprop through Cholesky)
- Primary: Solve for $\hat{\gamma}$ in terms of $\hat{F}_{\theta,\zeta}\varphi_{\zeta}(A,Z)$ and $\varphi_{\xi}(A)$
- Primary: Gradient steps on θ, ξ (backprop through Cholesky)
 - $\hat{F}_{\theta,\zeta}$ remains optimal wrt current φ_{θ} .

Iterate between updates of θ , ξ and ζ

Key point: features $\varphi_{\theta}(W)$ learned specially for:

$$\mathbb{E}(Y|a,z) = \mathbb{E}_{W|a,z} h_y(W,a)$$

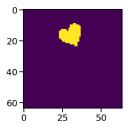
Contrast with autoencoders/sampling: must reconstruct/sample all of W.

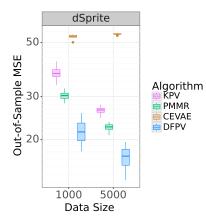
Experiments

Synthetic experiment, adaptive neural net features

dSprite example:

- $X = \{ scale, rotation, posX, posY \}$
- Treatment A is the image generated (with Gaussian noise)
- Outcome Y is quadratic function of A with multiplicative confounding by posY.
- Z = {scale, rotation, posX}, W = noisy image sharing posY
- Comparison with CEVAE (Louzios et al. 2017)



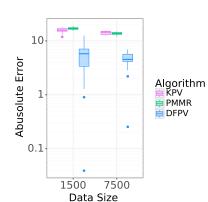


Louizos, Shalit, Mooij, Sontag, Zemel, Welling, Causal Effect Inference with Deep Latent-Variable 52/56 Models (2017)

Confounded offline policy evaluation

Synthetic dataset, demand prediction for flight purchase.

- Treatment A is ticket price.
- Policy $A \sim \pi(Z)$ depends on fuel price.



Conclusions

Neural net and kernel solutions:

- ...for ATE, CATE, dynamic treatment effects
- ...even for unobserved covariates/confounders (IV and proxy methods)
- ...with treatment A, covariates X, V, proxies/instruments (W, Z) multivariate, "complicated"
- Convergence guarantees for kernels and NN

Key messages:

- Don't meet your heroes model/sample hidden variables
- "Ridge regression is all you need"

Code available for all methods

Research support

Work supported by:

The Gatsby Charitable Foundation



Google Deepmind



Questions?

