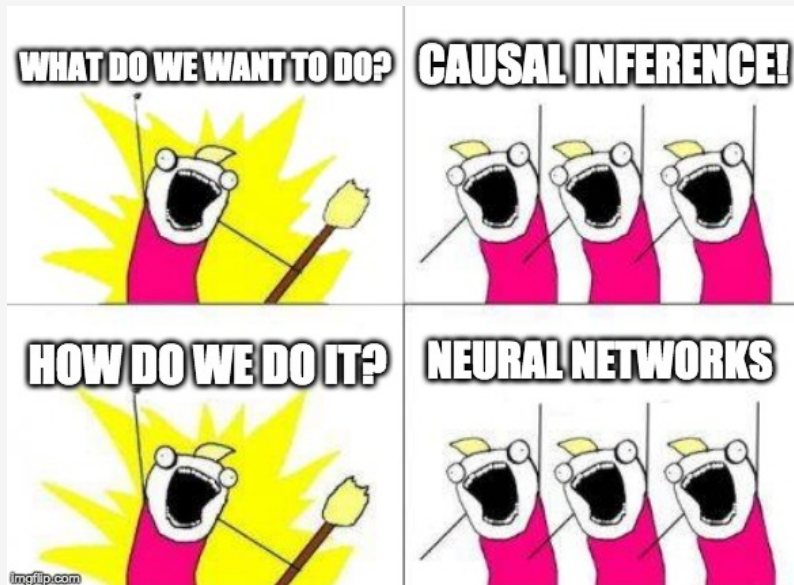


Adapting Neural Networks for the Estimation of Treatment Effects

Claudia Shi, David Blei, Victor Veitch

Columbia University

Causal Inference with Neural Nets using Observational Data



What is a Causal Question?

Questions about prediction:

- Will I have a headache tomorrow, given that I take this pill?
- What is the rate of drowning death, conditional on the ice cream sales is high?

Questions involve intervention:

- If I take this pill, will I have a headache tomorrow?
- Given that we increase the ice cream sales, what will the rate of drowning death be?

What is a Causal Question?

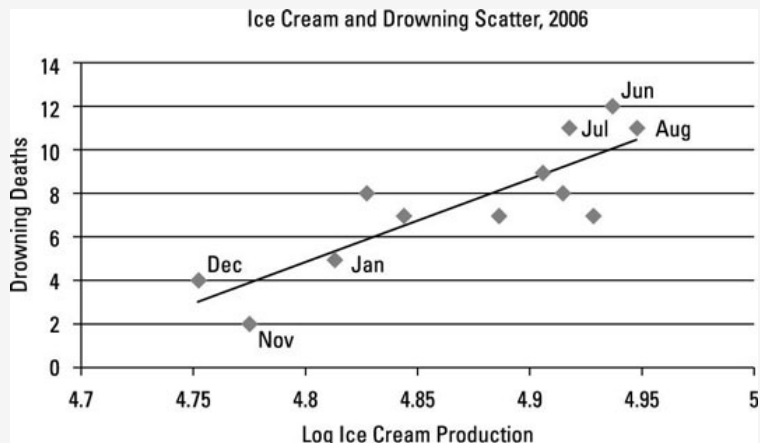
Questions about prediction:

- Will I have a headache tomorrow, given that I take this pill?
- What is the rate of drowning death, conditional on the ice cream sales is high?

Questions involve intervention:

- If I take this pill, will I have a headache tomorrow?
- Given that we increase the ice cream sales, what will the rate of drowning death be?

Observational Data: Confounding is a Problem



RCT Data: Not Accessible

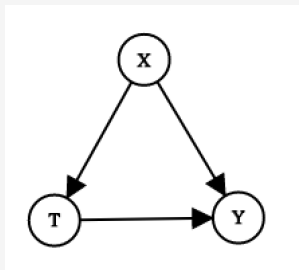


Causal Inference with Observational Data

Example

- treatment T (patient gets a drug)
- outcome Y (whether they recover)
- covariates X (illness severity, socioeconomic status)

What is expected effect of *intervening* by assigning the drug?

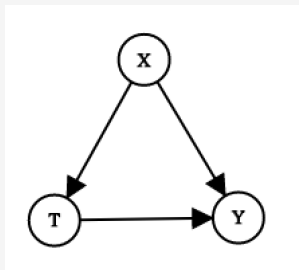


Causal Inference with Observational Data

Example

- treatment T (patient gets a drug)
- outcome Y (whether they recover)
- covariates X (illness severity, socioeconomic status)

What is expected effect of *intervening* by assigning the drug?



Adjustment

The average treatment effect is:

$$y = \mathbb{E}[Y \mid \text{do}(T = 1)] - \mathbb{E}[Y \mid \text{do}(T = 0)]$$

$$y \neq \mathbb{E}[Y \mid T = 1] - \mathbb{E}[Y \mid T = 0]$$

Theorem (No unobserved confounding)

If covariates X "block all backdoor paths" then

$$\begin{aligned} y &= \mathbb{E}[Y \mid \text{do}(T = 1)] - \mathbb{E}[Y \mid \text{do}(T = 0)] \\ &= \mathbb{E}[\mathbb{E}[Y \mid T = 1, X] - \mathbb{E}[Y \mid T = 0, X]] \end{aligned}$$

Adjustment

The average treatment effect is:

$$y = \mathbb{E}[Y \mid \text{do}(T = 1)] - \mathbb{E}[Y \mid \text{do}(T = 0)]$$

$$y \neq \mathbb{E}[Y \mid T = 1] - \mathbb{E}[Y \mid T = 0]$$

Theorem (No unobserved confounding)

If covariates X “block all backdoor paths” then

$$\begin{aligned} y &= \mathbb{E}[Y \mid \text{do}(T = 1)] - \mathbb{E}[Y \mid \text{do}(T = 0)] \\ &= \mathbb{E}[\mathbb{E}[Y \mid T = 1, X] - \mathbb{E}[Y \mid T = 0, X]] \end{aligned}$$

Average Treatment Effect:

$$y = \mathbb{E}[\mathbb{E}[Y | T = 1, X] - \mathbb{E}[Y | T = 0, X]]$$

Back Door Adjustment

Let $Q(t, x) = \mathbb{E}[Y | t, x]$, here is an estimator:

$$\hat{y}^Q = \frac{1}{n} \sum_i [\hat{Q}(1, x_i) - \hat{Q}(0, x_i)]$$

Estimation

Average Treatment Effect:

$$y = \mathbb{E}[\mathbb{E}[Y | T = 1, X] - \mathbb{E}[Y | T = 0, X]]$$

Back Door Adjustment

Let $Q(t, x) = \mathbb{E}[Y | t, x]$, here is an estimator:

$$\hat{y}^Q = \frac{1}{n} \sum_i [\hat{Q}(1, x_i) - \hat{Q}(0, x_i)]$$

Alternatively:

Average Treatment Effect:

$$y = \mathbb{E}[\mathbb{E}[Y | T = 1, X] - \mathbb{E}[Y | T = 0, X]]$$

Inverse Probability of Treatment Weighted Estimator (IPTW)

Let $g(x) = P(T = 1 | x)$, here is another estimator

$$\hat{y}^g = \frac{1}{n} \sum_i \left(\frac{t_i}{\hat{g}(x_i)} - \frac{1 - t_i}{1 - \hat{g}(x_i)} \right) y_i$$

Models

expected outcome: $Q(t, x) = \mathbb{E}[Y | t, x]$

propensity score: $g(x) = P(T = 1 | x)$

ATE: $y = \mathbb{E}[\mathbb{E}[Y | T = 1, X] - \mathbb{E}[Y | T = 0, X]]$

Semi-parametric efficient

- More complicated \hat{y} use both \hat{Q} and \hat{g}
- Nice asymptotic properties: low bias / efficient

Example: A Semi-parametric Efficient Estimator

Augmented IPTW

$$\hat{y} = \hat{Q}(1, x) - \hat{Q}(0, x) + \left(\frac{t}{\hat{g}(x)} - \frac{1-t}{1-\hat{g}(x)} \right) \{y - \hat{Q}(t, x)\}$$

This talk

We want to use neural networks to model Q and g .

How should we adapt the design and training of these networks so that \hat{y} is a good estimate of y ?

How should we adapt the design and training of these networks so that \hat{y} is a good estimate of y ?

Adaptations

- A neural network architecture—the Dragonnet—based on the sufficiency of the propensity score for causal estimation.
- A regularization procedure—targeted regularization—based on non-parametric estimation theory.

How should we adapt the design and training of these networks so that \hat{y} is a good estimate of y ?

Adaptations

- 1 A neural network architecture—the [Dragonnet](#)—based on the sufficiency of the propensity score for causal estimation.
- 2 A regularization procedure—[targeted regularization](#)—based on non-parametric estimation theory.

Dragonnet

Highlight

If the average treatment effect \boldsymbol{y} is identifiable from observational data by adjusting for X , then adjusting for the propensity score also succeeds.

Theorem (Rosenbaum and Rubin 1983)

If $\boldsymbol{y} = \mathbb{E}[\mathbb{E}[Y | T = 1, X] - \mathbb{E}[Y | T = 0, X]]$, then

$$\boldsymbol{y} = \mathbb{E}[\mathbb{E}[Y | T = 1, g(X)] - \mathbb{E}[Y | T = 0, g(X)]]$$

\implies estimate $\hat{Q}(t, x)$ using only parts of X relevant for T

Highlight

If the average treatment effect γ is identifiable from observational data by adjusting for X , then adjusting for the propensity score also succeeds.

Theorem (Rosenbaum and Rubin 1983)

If $\gamma = \mathbb{E}[\mathbb{E}[Y | T = 1, X] - \mathbb{E}[Y | T = 0, X]]$, then

$$\gamma = \mathbb{E}[\mathbb{E}[Y | T = 1, g(X)] - \mathbb{E}[Y | T = 0, g(X)]]$$

\implies estimate $\hat{Q}(t, x)$ using only parts of X relevant for T

One Natural Approach: Nednet

Goal

Estimate $\hat{Q}(t, \mathbf{x})$ using only parts of \mathbf{X} relevant for T

One Natural Approach: Nednet

Goal

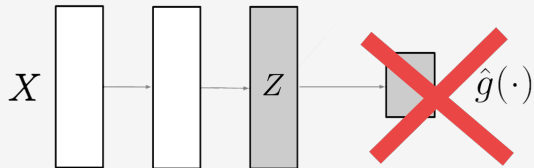
Estimate $\hat{Q}(t, \mathbf{x})$ using only parts of \mathbf{X} relevant for T



One Natural Approach: Nednet

Goal

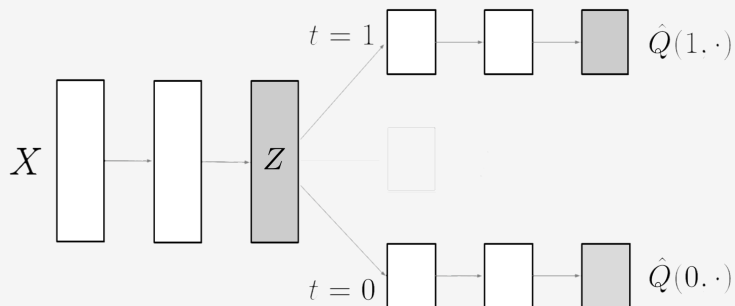
Estimate $\hat{Q}(t, \mathbf{x})$ using only parts of \mathbf{X} relevant for T



One Natural Approach: Nednet

Goal

Estimate $\hat{Q}(t, \mathbf{x})$ using only parts of \mathbf{X} relevant for T



One Natural Approach: Nednet

Goal

Estimate $\hat{Q}(t, x)$ using only parts of X relevant for T



Better: Dragonnet

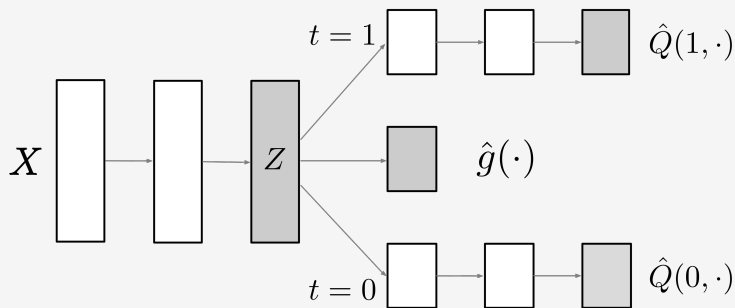
Goal: Estimate $\hat{Q}(t, \mathbf{x})$ using only parts of \mathbf{X} relevant for T

- Downstream estimator: $\hat{y}^Q = \frac{1}{n} \sum_i \left[\hat{Q}(1, \mathbf{x}_i) - \hat{Q}(0, \mathbf{x}_i) \right]$

Better: Dragonnet

Goal: Estimate $\hat{Q}(t, x)$ using only parts of X relevant for T

- Downstream estimator: $\hat{y}^Q = \frac{1}{n} \sum_i [\hat{Q}(1, x_i) - \hat{Q}(0, x_i)]$



Dragonnet

Is the End-to-end model better?

Is the End-to-end model better?

Table 4: Dragonnet produces more accurate estimates than NEDnet, a multi-stage alternative. Table entries are mean absolute error over all datasets.

IHDP	$\hat{\psi}^Q$	$\hat{\psi}^{\text{TMLE}}$
Dragonnet	0.12 ± 0.00	0.12 ± 0.00
NEDnet	0.15 ± 0.01	0.12 ± 0.00

ACIC	$\hat{\psi}^Q$	$\hat{\psi}^{\text{TMLE}}$
Dragonnet	0.55	1.97
NEDnet	1.49	2.80

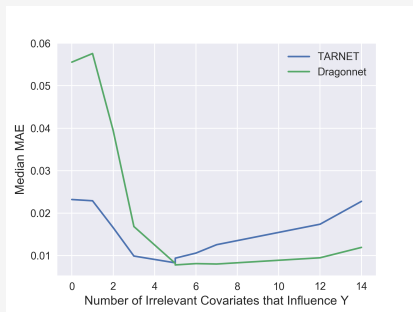
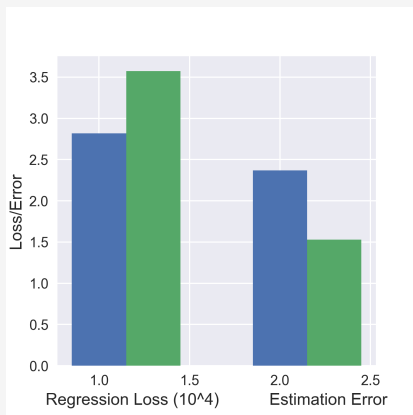
Does Dragonnet Actually Use Propensity Score Sufficiency?

- TARNET¹ = Dragonnet without treatment head
- $\hat{y}^Q = \frac{1}{n} \sum_i \left[\hat{Q}(1, x_i) - \hat{Q}(0, x_i) \right]$

¹<https://arxiv.org/abs/1606.03976>

Does Dragonnet Actually Use Propensity Score Sufficiency?

- TARNET¹ = Dragonnet without treatment head
- $\hat{y}^Q = \frac{1}{n} \sum_i [\hat{Q}(1, x_i) - \hat{Q}(0, x_i)]$



¹<https://arxiv.org/abs/1606.03976>

Targeted Regularization

Recall:

expected outcome: $Q(t, x) = \mathbb{E}[Y | t, x]$

propensity score: $g(x) = P(T = 1 | x)$

Average treatment effect: $\tau = \mathbb{E}[\mathbb{E}[Y | T = 1, X] - \mathbb{E}[Y | T = 0, X]]$

Semi-parametric efficient

- More complicated \hat{y} use both \hat{Q} and \hat{g}
- Nice asymptotic properties: low bias / efficient

Recall:

expected outcome: $Q(t, x) = \mathbb{E}[Y | t, x]$

propensity score: $g(x) = P(T = 1 | x)$

Average treatment effect: $\tau = \mathbb{E}[\mathbb{E}[Y | T = 1, X] - \mathbb{E}[Y | T = 0, X]]$

Semi-parametric efficient

- More complicated \hat{y} use both \hat{Q} and \hat{g}
- Nice asymptotic properties: low bias / efficient

Targeted Regularization

- targeted regularization is a modification to the objective function used for neural network training.
- based on non-parametric estimation theory.
- yields a fitted model that, with a suitable downstream estimator, guarantees desirable asymptotic properties.

Targeted Regularization

- targeted regularization is a modification to the objective function used for neural network training.
- based on non-parametric estimation theory.
- yields a fitted model that, with a suitable downstream estimator, guarantees desirable asymptotic properties.

Targeted Regularization

- targeted regularization is a modification to the objective function used for neural network training.
- based on non-parametric estimation theory.
- yields a fitted model that, with a suitable downstream estimator, guarantees desirable asymptotic properties.

Targeted Regularization

- targeted regularization is a modification to the objective function used for neural network training.
- based on **non-parametric estimation theory**.
- yields a fitted model that, with a suitable downstream estimator, guarantees desirable asymptotic properties.

Asymptotics

If $(\hat{Q}, \hat{g}, \hat{y})$ satisfy a certain equation, then

- **robustness** in the double machine-learning sense--- \hat{y} converges to y at a fast rate even if \hat{Q} and \hat{g} converge slowly
- **efficiency**---asymptotically, \hat{y} has the lowest variance of any consistent estimator of y

Asymptotics

If $(\hat{Q}, \hat{g}, \hat{y})$ satisfy a certain equation, then

- **robustness** in the double machine-learning sense--- \hat{y} converges to y at a fast rate even if \hat{Q} and \hat{g} converge slowly
- **efficiency**---asymptotically, \hat{y} has the lowest variance of any consistent estimator of y

Non-parametric Estimating Equation

Asymptotics hold if

- 1 \hat{Q} and \hat{g} are consistent
- 2 $(\hat{Q}, \hat{g}, \hat{y})$ satisfy non-parametric estimating equation,

$$0 = \frac{1}{n} \sum_i \dot{a} j (y_i, t_i, x_i; \hat{Q}, \hat{g}, \hat{y}),$$

where

$$\begin{aligned} j (y, t, x; Q, g, y) &= Q(1, x) - Q(0, x) \\ &+ \left(\frac{t}{g(x)} - \frac{1-t}{1-g(x)} \right) \{y - Q(t, x)\} - y \end{aligned}$$

A-IPTW

- 1 Obvious: fit \hat{Q} and \hat{g} , then choose \hat{y} so non-parametric estimating equation is satisfied

2

$$\hat{y} = \hat{Q}(1, x) - \hat{Q}(0, x) + \left(\frac{t}{\hat{g}(x)} - \frac{1-t}{1-\hat{g}(x)} \right) \{y - \hat{Q}(t, x)\}$$

- 3 Bad in practice; \hat{y} has $1/\hat{g}$ terms \implies finite sample = :(

A-IPTW

1 Obvious: fit \hat{Q} and \hat{g} , then choose \hat{y} so non-parametric estimating equation is satisfied

2

$$\hat{y} = \hat{Q}(1, x) - \hat{Q}(0, x) + \left(\frac{t}{\hat{g}(x)} - \frac{1-t}{1-\hat{g}(x)} \right) \{y - \hat{Q}(t, x)\}$$

3 Bad in practice; \hat{y} has $1/\hat{g}$ terms \implies finite sample = :(

A-IPTW

1 Obvious: fit \hat{Q} and \hat{g} , then choose \hat{y} so non-parametric estimating equation is satisfied

2

$$\hat{y} = \hat{Q}(1, x) - \hat{Q}(0, x) + \left(\frac{t}{\hat{g}(x)} - \frac{1-t}{1-\hat{g}(x)} \right) \{y - \hat{Q}(t, x)\}$$

3 Bad in practice; \hat{y} has $1/\hat{g}$ terms \implies finite sample = :(

A-IPTW

- 1 Obvious: fit \hat{Q} and \hat{g} , then choose \hat{y} so non-parametric estimating equation is satisfied

- 2

$$\hat{y} = \hat{Q}(1, \mathbf{x}) - \hat{Q}(0, \mathbf{x}) + \left(\frac{t}{\hat{g}(\mathbf{x})} - \frac{1-t}{1-\hat{g}(\mathbf{x})} \right) \{y - \hat{Q}(t, \mathbf{x})\}$$

- 3 Bad in practice; \hat{y} has $1/\hat{g}$ terms \implies finite sample = :(

Alternative

- Choose $\hat{y}^Q = \frac{1}{n} \sum_i [\hat{Q}(1, x_i) - \hat{Q}(0, x_i)]$ —no bad $1/\hat{g}$ terms
- Fit \hat{Q} and \hat{g} so non-parametric estimating equation is satisfied

A-IPTW

- 1 Obvious: fit \hat{Q} and \hat{g} , then choose \hat{y} so non-parametric estimating equation is satisfied

2

$$\hat{y} = \hat{Q}(1, \mathbf{x}) - \hat{Q}(0, \mathbf{x}) + \left(\frac{t}{\hat{g}(\mathbf{x})} - \frac{1-t}{1-\hat{g}(\mathbf{x})} \right) \{y - \hat{Q}(t, \mathbf{x})\}$$

- 3 Bad in practice; \hat{y} has $1/\hat{g}$ terms \implies finite sample = :(

Alternative

- Choose $\hat{y}^Q = \frac{1}{n} \sum_i \left[\hat{Q}(1, \mathbf{x}_i) - \hat{Q}(0, \mathbf{x}_i) \right]$ —no bad $1/\hat{g}$ terms
- Fit \hat{Q} and \hat{g} so non-parametric estimating equation is satisfied

Targeted Regularization

Introduce extra parameter \mathbf{e} and regularization term $g(y, t, \mathbf{x}; \mathbf{q}, \mathbf{e})$

$$\tilde{Q}(t_i, \mathbf{x}_i; \mathbf{q}, \mathbf{e}) = Q^{\text{nn}}(t_i, \mathbf{x}_i; \mathbf{q}) + \mathbf{e} \left[\frac{t_i}{g^{\text{nn}}(\mathbf{x}_i; \mathbf{q})} - \frac{1 - t_i}{1 - g^{\text{nn}}(\mathbf{x}_i; \mathbf{q})} \right]$$

$$g(y_i, t_i, \mathbf{x}_i; \mathbf{q}, \mathbf{e}) = (y_i - \tilde{Q}(t_i, \mathbf{x}_i; \mathbf{q}, \mathbf{e}))^2$$

Then train as

$$\hat{\mathbf{q}}, \hat{\mathbf{e}} = \underset{\mathbf{q}, \mathbf{e}}{\operatorname{argmin}} \left[\underbrace{\hat{R}(\mathbf{q}; \mathbf{X})}_{\text{usual objective}} + a \underbrace{\frac{1}{n} \sum_i g(y_i, t_i, \mathbf{x}_i; \mathbf{q}, \mathbf{e})}_{\text{targeted regularization}} \right]$$

Define an estimator \hat{y}^{treg} ,

$$\hat{y}^{\text{treg}} = \frac{1}{n} \mathring{a}_i \hat{Q}^{\text{treg}}(1, x_i) - \hat{Q}^{\text{treg}}(0, x_i), \quad \text{where}$$

$$\hat{Q}^{\text{treg}} = \tilde{Q}(\cdot, \cdot; \hat{q}, \hat{e}).$$

Define an estimator \hat{y}^{treg} ,

$$\hat{y}^{\text{treg}} = \frac{1}{n} \hat{a}_i \hat{Q}^{\text{treg}}(1, x_i) - \hat{Q}^{\text{treg}}(0, x_i), \quad \text{where}$$

$$\hat{Q}^{\text{treg}} = \tilde{Q}(\cdot, \cdot; \hat{q}, \hat{e}).$$

The point is:

$$0 = \mathbb{1}_e(\hat{R}(q; \mathbf{X}) + a \frac{1}{n} \hat{a}_i g(y_i, t_i, x_i; q, e))|_{\hat{e}} = a \frac{1}{n} \hat{a}_i j(y_i, t_i, x_i; \hat{Q}^{\text{treg}}, \hat{g}, \hat{y}^{\text{treg}}).$$

Define an estimator \hat{y}^{treg} ,

$$\hat{y}^{\text{treg}} = \frac{1}{n} \hat{a}_i \hat{Q}^{\text{treg}}(1, x_i) - \hat{Q}^{\text{treg}}(0, x_i), \quad \text{where}$$

$$\hat{Q}^{\text{treg}} = \tilde{Q}(\cdot, \cdot; \hat{q}, \hat{e}).$$

The point is:

$$0 = \mathbb{1}_e(\hat{R}(q; \mathbf{X}) + a \frac{1}{n} \hat{a}_i g(y_i, t_i, x_i; q, e))|_{\hat{e}} = a \frac{1}{n} \hat{a}_i j(y_i, t_i, x_i; \hat{Q}^{\text{treg}}, \hat{g}, \hat{y}^{\text{treg}}).$$

That is, minimizing the targeted regularization term forces $(\hat{Q}^{\text{treg}}, \hat{g}, \hat{y}^{\text{treg}})$ to satisfy the non-parametric estimating equation.

Experiment

Infant Health Development Program Benchmark (IHDP)

Method	Δ_{in}	Δ_{out}	Δ_{all}
BNN [JSS16]	$0.37 \pm .03$	$0.42 \pm .03$	—
TARNET [SJS16]	$0.26 \pm .01$	$0.28 \pm .01$	—
CFR Wass[SJS16]	$0.25 \pm .01$	$0.27 \pm .01$	—
CEVAEs [Lou+17]	$0.34 \pm .01$	$0.46 \pm .02$	—
GANITE [YJS18]	$0.43 \pm .05$	$0.49 \pm .05$	—
baseline (TARNET)	$0.16 \pm .01$	$0.21 \pm .01$	$0.13 \pm .00$
baseline + t-reg	$0.15 \pm .01$	$0.20 \pm .01$	$0.12 \pm .00$
Dragonnet	$0.14 \pm .01$	$0.21 \pm .01$	$0.12 \pm .00$
Dragonnet + t-reg	$0.14 \pm .01$	$0.20 \pm .01$	$0.11 \pm .00$

IBM Causal Inference Benchmarking Framework (ACIC)

Table 2: Dragonnet and targeted regularization improve estimation on average on ACIC 2018. Table entries are mean absolute error over all datasets.

Method	Δ_{all}
baseline (TARNET)	1.45
baseline + t-reg	1.40
Dragonnet	0.55
Dragonnet + t-reg	0.35

Table 3: Dragonnet and targeted regularization improve over the baseline about half the time, but improvement is substantial when it does happen. Error values are mean absolute error on ACIC 2018.

ψ^Q	% <i>improve</i>	\uparrow_{avg}	\downarrow_{avg}
baseline:	0%	0	0
+ t-reg	42%	0.30	0.11
+ dragon	63%	1.42	0.01
+ dragon & t-reg	46%	2.37	0.01

Summary

- Dragonnet: a neural network architecture based on the sufficiency of the propensity score for causal estimation.
- targeted regularization: a regularization procedure based on non-parametric estimation theory.
- They both work!

Thank You!

- Adapting Neural Networks for the Estimation of Treatment Effects.
[arxiv:1906.02120](https://arxiv.org/abs/1906.02120)