# CSC2541: Introduction to Causality Lecture 7 - Double Machine Learning

Instructor: Rahul G. Krishnan TA: Vahid Balazadeh-Meresht

November 21, 2022

### Back to ML for Causality

- ▶ Last lecture: Use-case of the invariance assumption in causal inference for machine learning (ML)
- ▶ Today: How to use ML predictive models to get *unbiased* causal effect estimations with *fast convergence rate* and *confidence intervals*?
- ▶ TAR-Net also used ML for causal effect estimation. However:
  - Not flexible in using different ML models,
  - No convergence rate guarantees,
  - No uncertainty regions,
  - Only for binary treatments.
- We will assume ignorability, i.e., covariates X block all the backdoor paths from treatment T to outcome Y

# Where can we use ML for causal estimation?

- ▶ ML methods are effective in prediction contexts, but this does not translate into good performance for estimation of "causal" parameters
  - 1. Overfitting bias: Capturing more than the relationship of T and Y
  - 2. Regularization bias: Slower convergence rate
- $\blacktriangleright$  Often, covariates X are high-dimensional while T is low-dimensional
- ▶ The relationship between Y and X is more complex than the relationship between Y and T
- $\blacktriangleright$  Idea: Use ML methods to model  $Y \sim X$  and linear models for  $Y \sim T$

## A canonical example - Partially Linear Model

Assume the following data generating process:

$$\begin{split} Y &= \alpha_0 T + g_0(X) + U \\ T &= m_0(X) + V \\ \text{with } \mathbb{E}\left[U|T,X\right] = 0, \ \mathbb{E}\left[V|X\right] = 0 \end{split}$$

 $\blacktriangleright Y, T \in \mathbb{R}$ 

- $\alpha_0$  is the target parameter of interest (ATE)
- $\blacktriangleright$  X is a high-dimensional vector
- We call  $\eta_0 = (g_0, m_0)$  nuisance parameters We do not care about their estimation as long as it results in correct  $\alpha_0$

## Naive prediction-based ML approach is Bad

 $\blacktriangleright$  Predict Y using X and T:

$$\hat{Y} = \hat{\alpha}_0 T + \hat{g}_0(X)$$

▶ For example, we can fit the model by alternating minimization

- Given initial parameters, run a Random Forest on  $Y \hat{\alpha}_0 T$  to fit  $\hat{g}_0(X)$
- ▶ Run Ordinary Least Squares (OLS) on  $Y \hat{g}_0(X)$  to fit  $\hat{\alpha}_0$
- Repeat until convergence
- ► Good prediction performance  $\|\hat{Y} Y\|_2^2$ . But, the distribution of  $\alpha_0 \hat{\alpha}_0$  looks like this



## Why is the naive approach bad?

- Assume the minimization is converged and we learned  $\hat{g}_0(X)$
- $\hat{\alpha}_0$  is the OLS solution to  $Y = \alpha T + \hat{g}_0(X)$ :

$$\hat{\alpha}_0 = \left(\frac{1}{n}\sum_i T_i^2\right)^{-1} \frac{1}{n}\sum_i T_i(Y_i - \hat{g}_0(X_i))$$
  
assuming  $\mathbb{E}[a_0(X)] = \mathbb{E}[m_0(X)] = 0$ 

assuming 
$$\mathbb{E}[g_0(X)] = \mathbb{E}[m_0(X)] = 0$$

► Let's look at the error:  

$$\hat{\alpha}_{0} = \left(\frac{1}{n}\sum_{i}T_{i}^{2}\right)^{-1}\frac{1}{n}\sum_{i}T_{i}(Y_{i} - \hat{g}_{0}(X_{i}))$$

$$= \left(\frac{1}{n}\sum_{i}T_{i}^{2}\right)^{-1}\frac{1}{n}\sum_{i}T_{i}(\alpha_{0}T_{i} + g_{0}(X_{i}) + U_{i} - \hat{g}_{0}(X_{i}))$$

$$= \left(\frac{1}{n}\sum_{i}T_{i}^{2}\right)^{-1}\left[\left(\frac{1}{n}\sum_{i}T_{i}^{2}\right)\alpha_{0} + \left(\frac{1}{n}\sum_{i}T_{i}U_{i}\right) + \left(\frac{1}{n}\sum_{i}T_{i}(g_{0}(X_{i}) - \hat{g}_{0}(X_{i})\right)\right]\right]$$

$$= \alpha_{0} + \left(\frac{1}{n}\sum_{i}T_{i}^{2}\right)^{-1}\left[\frac{1}{n}\sum_{i}T_{i}U_{i} + \frac{1}{n}\sum_{i}T_{i}(g_{0}(X_{i}) - \hat{g}_{0}(X_{i}))\right]$$

$$= \alpha_{0} + \left(\frac{1}{n}\sum_{i}T_{i}^{2}\right)^{-1}\left[\frac{1}{n}\sum_{i}T_{i}U_{i} + \frac{1}{n}\sum_{i}\left(m_{0}(X_{i}) + V_{i}\right)(g_{0}(X_{i}) - \hat{g}_{0}(X_{i})\right]$$

$$= \left(\frac{1}{n}\sum_{i}\left[T_{i}^{2}\right]^{-1}\right]^{-1}\left[\frac{1}{n}\sum_{i}T_{i}U_{i} + \frac{1}{n}\sum_{i}\left(m_{0}(X_{i}) + V_{i}\right)(g_{0}(X_{i}) - \hat{g}_{0}(X_{i})\right]$$

## Why is the naive approach bad?

$$= \underbrace{\begin{bmatrix} A \\ \hline \frac{1}{\sqrt{n}\sum_{i}T_{i}U_{i}} + \frac{1}{\sqrt{n}\sum_{i}m_{0}(X_{i})(g_{0}(X_{i}) - \hat{g}_{0}(X_{i}))} + \frac{1}{\sqrt{n}\sum_{i}V_{i}(g_{0}(X_{i}) - \hat{g}_{0}(X_{i}))} \end{bmatrix}}_{\mathbb{E}\left[T_{i}^{2}\right]}$$

The goal is to find a root-n consistent and asymptotically normal estimate of  $\alpha_0$ , i.e.,  $\sqrt{n}(\hat{\alpha}_0 - \alpha_0) \rightarrow \mathcal{N}(0, \sigma^2)$ 

- $A \to \mathcal{N}(0, \sigma_A^2)$  by Central Limit Theorem. It can be seen as sample average of random variables  $T_i U_i$
- What about term B? Does  $B \to \mathcal{N}(0, \sigma_B^2)$  for some  $\sigma_B^2$ ?

## Regularization Bias - Term B

Machine learning methods employ regularization (e.g., L<sup>2</sup> regularization) to reduce variance. However, this often induces bias and lower convergence rate:

$$g_0(X_i) - \hat{g}_0(X_i) \propto n^{-\phi_g}$$
, for some  $\phi_g < \frac{1}{2}$  (slow convergence)

Therefore, term B will be

$$B = \frac{1}{\sqrt{n}} \sum_{i} m_0(X_i) \left( g_0(X_i) - \hat{g}_0(X_i) \right) \propto \frac{1}{\sqrt{n}} \cdot n \cdot n^{-\phi_g} \propto n^{\frac{1}{2} - \phi_g} \to \infty$$

▶ How to make this term vanish?

### Double Machine Learning

- ► The naive approach was the OLS solution to  $Y = \alpha T + \hat{g}_0(X)$
- ▶ Idea: Partial out the effect of covariate X on treatment T
  - Train an ML algorithm to predict T from X:  $\hat{T} = \hat{m}_0(X)$
  - Consider the residual  $\hat{V} = T \hat{m}_0(X)$
  - Find the OLS solution  $\hat{\beta}$  to  $Y = \beta \hat{V} + \hat{g}_0(X)$
- ▶ This approach is called Double Machine Learning (DML) as we use machine learning twice: to learn  $\hat{g}_0(X)$  and to learn  $\hat{m}_0(X)$
- ▶  $\hat{\beta}$  is a root-n consistent estimate of  $\alpha_0$ .  $(\alpha_0 \hat{\beta})$  looks like this



### Partialling out the effect of covariates. Frisch-Waugh-Lovell theorem

- But why does partialling out the effect of X on T results in a valid estimate?
- ▶ Let's make everything linear. Consider the following linear equation:

$$Y = T\beta_1 + X\beta_2$$

for  $T, Y, \beta_1 \in \mathbb{R}$  and  $\beta_2, X \in \mathbb{R}^d$ . Assume Y, T, X are data matrices

- To estimate  $\beta_1$ , one can use OLS by concatenating T and X
- Frisch–Waugh–Lovell (FWL) theorem says we can estimate  $\beta_1$  in another way. Residuals-on-residuals:
  - Regress (linear)  $\boldsymbol{Y}$  on  $\boldsymbol{X}$  and let  $\hat{\boldsymbol{U}} = \boldsymbol{Y} \hat{\boldsymbol{Y}}$
  - Regress (linear) T on X and let  $\hat{V} = T \hat{T}$
  - Regress (linear)  $\hat{U}$  on  $\hat{V}$  to estimate  $\beta_1$
- ▶ FWL is a simpler version of DML. Instead of arbitrary ML methods, it uses linear regression

### FWL theorem - Proof

- Define the prediction matrix  $\boldsymbol{P} = \boldsymbol{X} (\boldsymbol{X}^{\top} \boldsymbol{X})^{-1} \boldsymbol{X}^{\top}$ 
  - ▶ E.g., the OLS solution for  $Y \sim X$ :  $\hat{Y} = X(X^{\top}X)^{-1}X^{\top}Y = PY$
- ▶ Define the residual matrix R = I P

• E.g., 
$$\boldsymbol{Y} - \hat{\boldsymbol{Y}} = \boldsymbol{Y} - \boldsymbol{P}\boldsymbol{Y} = \boldsymbol{R}\boldsymbol{Y}$$

▶ Note that residuals are **orthogonal** to predicted values

$$RP = (I - P)P = P - P^{2} = 0$$

• Let's apply the residual matrix on  $Y = T\beta_1 + X\beta_2$ :

$$RY = RT\beta_1 + RX\beta_2$$

► However,

$$\boldsymbol{R}\boldsymbol{X} = (\boldsymbol{I} - \boldsymbol{X}(\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top})\boldsymbol{X} = \boldsymbol{X} - \boldsymbol{X}(\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}\boldsymbol{X} = \boldsymbol{0}$$

► Therefore,

$$oldsymbol{R} oldsymbol{Y} = oldsymbol{R} oldsymbol{T} eta_1$$
  
 $oldsymbol{Y} - \hat{oldsymbol{Y}} = (oldsymbol{T} - \hat{oldsymbol{T}})eta_1$ 

# Back to DML - Overcoming regularization bias

- FWL shows that partialling out X does not affect the relationship between Y and T. It essentially gives the **same** answer
- But why does the estimation from DML  $(\hat{\beta})$  converges **better** than the naive solution  $\hat{\alpha}_0$ ?
- The key is the regularization bias (term B)

$$= \frac{\begin{bmatrix} A \\ \hline \frac{1}{\sqrt{n}\sum_{i}T_{i}U_{i}} + \frac{B}{\sqrt{n}\sum_{i}m_{0}(X_{i})(g_{0}(X_{i}) - \hat{g}_{0}(X_{i}))} + \frac{C}{\sqrt{n}\sum_{i}V_{i}(g_{0}(X_{i}) - \hat{g}_{0}(X_{i}))} \end{bmatrix}}{\mathbb{E}\left[T_{i}^{2}\right]}$$

• Let's write a similar estimation error for the DML solution  $\hat{\beta}$ 

# Why is the DML approach good?

▶  $\hat{\beta}_0$  is the OLS solution to  $Y = \beta \hat{V} + \hat{g}_0(X)$ , where  $\hat{V} = T - \hat{m}_0(X)$ 

$$\hat{\beta}_0 = \left(\frac{1}{n}\sum_i \hat{V}_i^2\right)^{-1} \frac{1}{n}\sum_i \hat{V}_i(Y_i - \hat{g}_0(X_i))$$

▶ For a simpler analysis, we consider a slightly different estimator

$$\hat{\beta} = \left(\frac{1}{n}\sum_{i}\hat{V}_{i}T_{i}\right)^{-1}\frac{1}{n}\sum_{i}\hat{V}_{i}(Y_{i}-\hat{g}_{0}(X_{i}))$$

▶ In finite samples,  $\hat{\beta} \neq \hat{\beta}_0$ . However, they both will have similar asymptotic properties as  $\mathbb{E}[\hat{V}^2] = \mathbb{E}[\hat{V}T]$  for infinite samples

# Why is the DML approach good?

Let's look at the error:

$$\begin{split} \hat{\beta} &= \left(\frac{1}{n}\sum_{i}\hat{V}_{i}T_{i}\right)^{-1}\frac{1}{n}\sum_{i}\hat{V}_{i}(Y_{i}-\hat{g}_{0}(X_{i}))\\ &= \left(\frac{1}{n}\sum_{i}\hat{V}_{i}T_{i}\right)^{-1}\frac{1}{n}\sum_{i}\hat{V}_{i}(\alpha_{0}T_{i}+g_{0}(X_{i})+U_{i}-\hat{g}_{0}(X_{i}))\\ &= \alpha_{0} + \left(\frac{1}{n}\sum_{i}\hat{V}_{i}T_{i}\right)^{-1}\left[\left(\frac{1}{n}\sum_{i}\hat{V}_{i}U_{i}\right) + \left(\frac{1}{n}\sum_{i}\hat{V}_{i}(g_{0}(X_{i})-\hat{g}_{0}(X_{i})\right)\right]\right]\\ &= \alpha_{0} + \left(\frac{1}{n}\sum_{i}\hat{V}_{i}T_{i}\right)^{-1}\left[\left(\frac{1}{n}\sum_{i}\hat{V}_{i}U_{i}\right) + \left(\frac{1}{n}\sum_{i}(T_{i}-\hat{m}_{0}(X_{i}))(g_{0}(X_{i})-\hat{g}_{0}(X_{i})\right)\right]\right]\\ &= \alpha_{0} + \frac{\left[\left(\frac{1}{n}\sum_{i}\hat{V}_{i}U_{i}\right) + \left(\frac{1}{n}\sum_{i}(m_{0}(X_{i})+V_{i}-\hat{m}_{0}(X_{i}))(g_{0}(X_{i})-\hat{g}_{0}(X_{i}))\right)\right]}{\left(\frac{1}{n}\sum_{i}\hat{V}_{i}T_{i}\right)} \end{split}$$

# Why is the DML approach good?

▶ Therefore,

$$\boxed{ \underbrace{ \begin{bmatrix} A' \\ \hline \frac{1}{\sqrt{n} \sum_{i} \hat{V}_{i} U_{i}} + \underbrace{\frac{1}{\sqrt{n} \sum_{i} (m_{0}(X_{i}) - \hat{m}_{0}(X_{i})) (g_{0}(X_{i}) - \hat{g}_{0}(X_{i}))}_{\left(\frac{1}{n} \sum_{i} \hat{V}_{i} (g_{0}(X_{i}) - \hat{g}_{0}(X_{i}))\right)} } \underbrace{ \begin{bmatrix} \frac{1}{n} \sum_{i} \hat{V}_{i} (g_{0}(X_{i}) - \hat{g}_{0}(X_{i})) \\ \hline \frac{1}{\sqrt{n} \sum_{i} \hat{V}_{i} T_{i}} \end{bmatrix} }$$

$$\begin{bmatrix} A \\ \hline \frac{1}{\sqrt{n}\sum_{i}^{n}T_{i}U_{i}} + \frac{1}{\sqrt{n}\sum_{i}^{n}m_{0}(X_{i})(g_{0}(X_{i}) - \hat{g}_{0}(X_{i}))} + \frac{1}{\sqrt{n}\sum_{i}^{n}V_{i}(g_{0}(X_{i}) - \hat{g}_{0}(X_{i}))} \end{bmatrix}$$
$$\mathbb{E}[T_{i}^{2}]$$

• A' behaves similarly to A. Term C is exactly the same. The difference is in the regularization terms B and B'

### DML overcomes the regularization bias

$$B' = \frac{1}{\sqrt{n}} \sum_{i} (m_0(X_i) - \hat{m}_0(X_i)) (g_0(X_i) - \hat{g}_0(X_i))$$

• Again, since we are using (regularized) ML methods, the convergence rates of  $g_0$  and  $m_0$  are slow

$$g_0(X_i) - \hat{g}_0(X_i) \propto n^{-\phi_g}, \text{ for some } \phi_g < \frac{1}{2}$$
$$m_0(X_i) - \hat{m}_0(X_i) \propto n^{-\phi_m}, \text{ for some } \phi_m < \frac{1}{2}$$

ы

Therefore, term B' will be

$$B' = \frac{1}{\sqrt{n}} \sum_{i} (m_0(X_i) - \hat{m}_0(X_i)) (g_0(X_i) - \hat{g}_0(X_i))$$
$$\propto \frac{1}{\sqrt{n}} \cdot n \cdot n^{-\phi_m} \cdot n^{-\phi_g} \propto n^{\frac{1}{2} - \phi_g - \phi_m}$$

▶ Now, even for slow convergence rates like  $\phi_g, \phi_m = \frac{1}{4} + \epsilon, B'$  will converge with root-n rate

$$n^{\frac{1}{2}-\phi_g-\phi_m} = n^{\frac{1}{2}-\frac{1}{4}-\epsilon-\frac{1}{4}-\epsilon} = n^{-2\epsilon} \to 0$$

## Overfitting bias - Term ${\cal C}$

$$\begin{split} \sqrt{n}(\hat{\beta} - \alpha_0) &= \\ \underbrace{\begin{bmatrix} A' \\ \frac{1}{\sqrt{n}\sum_{i} \hat{V}_i U_i} + \frac{1}{\sqrt{n}\sum_{i} (m_0(X_i) - \hat{m}_0(X_i)) \left(g_0(X_i) - \hat{g}_0(X_i)\right)} + \frac{C}{\sqrt{n}\sum_{i} V_i \left(g_0(X_i) - \hat{g}_0(X_i)\right)} \\ \frac{(\frac{1}{n}\sum_{i} \hat{V}_i T_i)} \end{split} \end{split}$$

• We saw 
$$A' \to \mathcal{N}(0, \sigma_A^2)$$

▶ DML used orthogonalization to overcome regularization bias B':  $B' \to \mathcal{N}(0, \sigma_B^2)$ 

• What about term C? Does it also vanish?

### Overfitting bias - Term C

- ▶ To learn  $\hat{g}_0(X)$ , we fitted an ML method to predict Y from X
- ▶ For example, we can (artificially) assume the estimator is as follows

$$\hat{g}_0(X_i) = g_0(X_i) + \underbrace{\frac{(Y_i - g_0(X_i))}{n^{1/2 - \epsilon}}}_{\text{error}} \qquad \text{(fast but not root-n rate)}$$

The error term is the part of Y that is unexplainable by g<sub>0</sub>(X)
 Let's look at term C:

$$\begin{split} C &= \frac{1}{\sqrt{n}} \sum_{i} V_{i} \left( g_{0}(X_{i}) - \hat{g}_{0}(X_{i}) \right) \\ &= \frac{1}{\sqrt{n}} \sum_{i} V_{i} \frac{(Y_{i} - g_{0}(X_{i}))}{n^{1/2 - \epsilon}} \\ &= \frac{1}{\sqrt{n}} \sum_{i} V_{i} \frac{(g_{0}(X_{i}) + T_{i} + U_{i} - g_{0}(X_{i}))}{n^{1/2 - \epsilon}} \\ &= \frac{1}{\sqrt{n}} \sum_{i} V_{i} \frac{(T_{i} + U_{i})}{n^{1/2 - \epsilon}} \\ &= \frac{1}{\sqrt{n}} \sum_{i} V_{i} \frac{(m_{0}(X_{i}) + V_{i} + U_{i})}{n^{1/2 - \epsilon}} = \frac{1}{\sqrt{n}} \sum_{i} \frac{V_{i}^{2}}{n^{1/2 - \epsilon}} + \ldots = \frac{n}{\sqrt{n}n^{1/2 - \epsilon}} \sum_{i} \frac{V_{i}^{2}}{n} + \ldots \\ &= n^{\epsilon} \operatorname{Var}(V) + \ldots \\ &\to \infty \end{split}$$

### Removing the overfitting bias with Sample Splitting

- ▶ Term C explodes since the estimated  $\hat{g}_0(X)$  is overfitted: It captures more than  $g_0(X)$  from Y and becomes related to noise V
- ▶ To overcome this, DML uses sample splitting
  - Use part of samples  $(I \subset \{1, 2, ..., n\})$  to estimate  $\hat{\beta}$
  - Use auxiliary samples  $(I^c)$  to estimate  $\hat{g}_0(X)$
- $\blacktriangleright$  Therefore, term C will be

$$C = \frac{1}{\sqrt{n}} \sum_{i \in I} V_i(g_0(X_i) - \hat{g}_0(X_i))$$

 $\blacktriangleright$  This new C will vanish. Let's look at it's expectation

$$\mathbb{E}[C] = \frac{1}{\sqrt{n}} \sum_{i \in I} \mathbb{E}\left[V_i(g_0(X_i) - \hat{g}_0(X_i))\right]$$
  
$$= \frac{1}{\sqrt{n}} \sum_{i \in I} \mathbb{E}[\mathbb{E}[V_i(\underline{g_0(X_i) - \hat{g}_0(X_i)})] |X_{I^c}]] \quad \text{(condition on auxiliary samples)}$$
  
$$= \frac{1}{\sqrt{n}} \sum_{i \in I} \mathbb{E}[\mathbb{E}[V_i]\mathbb{E}[Err_i|X_{I^c}]] \quad (Err_i \text{ only depends on auxiliary samples)}$$
  
$$= 0 \qquad (\mathbb{E}[V_i] = 0)$$

### DML Algorithm - Summary

In summary, for a given dataset  $\{T^i, X^i, Y^i\}_{i=1}^n$ , DML follows the following to estimate average treatment effect:

- 1. Split samples to two parts I and  $I^c$  s.t.  $I\cup I^c=\{1,\ldots,n\}$  and  $I\cap I^c=\emptyset$
- 2. Train any (regularized) machine learning model  $M_t$  to predict T from X using auxiliary  $I^c$
- 3. Train any (regularized) machine learning model  $M_y$  to predict Y from X using  $I^c$
- 4. Obtain the residuals  $Y_R = Y M_y(X)$  and  $T_R = T M_t(X)$  from samples I
- 5. Regress (linearly)  $Y_R$  on  $T_R$  to get the estimated ATE

To increase sample efficiency, we can get another estimate by changing the role of I and  $I^c$  and take the average of the two estimations

## DML properties

- ► It allows using any a broad range of ML or non-parametric algorithms to estimate high-dimensional nuisance parameters  $(\eta_0 = (g_0, m_0))$
- ▶ It gives a root-n consistent estimator for ATE Fast convergence
- ▶ We can get valid confidence intervals over ATE as the estimate is asymptotically normal
- $\blacktriangleright$  DML was published in 2016<sup>1</sup> and is still among the best methods in causal inference competitions<sup>2</sup>

<sup>&</sup>lt;sup>1</sup>Chernozhukov et al., 2016.

<sup>&</sup>lt;sup>2</sup>ACIC 2022 data challenge - https://acic2022.mathematica.org/results

## Application outside of causal estimation - Identifying causal parents

- Now that we know what double machine learning actually is, how can we use it to solve practical problems?
- **Question:** How do we identify the causal parents of a variable?
- Given genetic expression data, we might want to know which are the causal parents while controlling for the effect of other genes.
- (Raj et al., 2020) use DML as a black box and devise a parallel search strategy to treat each gene as a treatment and predict the outcome (disease incidence).
- **Discuss:** What might the pros/cons of this approach be?

## Recap of Introduction to Causality

- Correlation is not causation!
- No causal inference without assumptions positivity, no unobserved confounding.
- Potential outcomes, Causal Bayesian networks, Structural causal models.
- ▶ Identifying interventions, evaluating counterfactuals and do-calculus.
- Estimation methods: G-formula, Matching, Inverse propensity weighting.
- Handling unobserved confounding instrumental variables and local average treatment effects.
- ▶ Causal inference for ML: learning from environments.
- ML for causal inference: double machine learning for estimating treatment effects.

### What we did not cover

- Sensitivity analysis understanding how much unobserved confounding one needs to change the outcomes of your study.
- ▶ Dynamic treatment effects causal effects with time-varying data.
- Partial identification bounding causal effects rather than point identification.
- Causal decision making what (among) many interventions should I make?
- Applications of causal inference to improve RL, control, planning, predictive modeling in healthcare.
- ▶ Causal representation learning ????

### General advice



Figure: Be critical of the methods you use!

- The easiest person to fool is yourself always question your assumptions!
- Work closely with domain experts common sense and practical wisdom >>> any result from any algorithm.
- ▶ Always ask "where do the bits come from"?

Chernozhukov, Victor et al. (2016). "Double/debiased machine learning for treatment and causal parameters". In: *arXiv preprint arXiv:1608.00060*.

**R**aj, Anant et al. (2020). "Causal feature selection via orthogonal search". In: *arXiv preprint arXiv:2007.02938*.