



# Causal Machine Learning

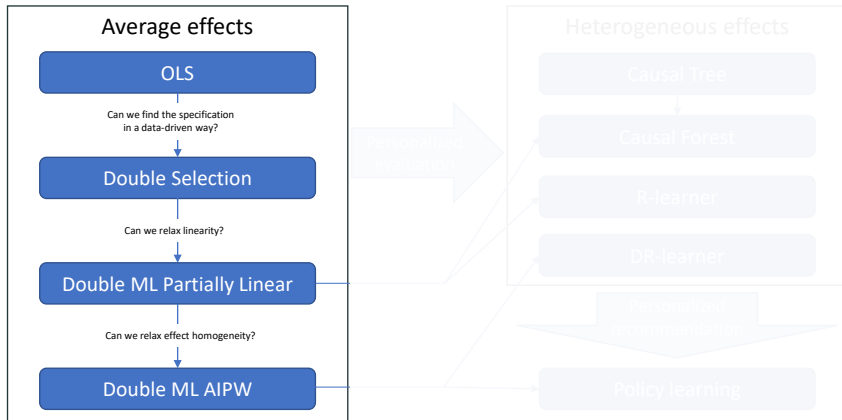
Double ML - the general recipe

---

Michael Knaus

March 13, 2023

# Current state of affairs



Double ML is a generic recipe and can be used for other target parameters/research designs

# Plan of this morning

Understand the generic recipe of Double ML

1. The Double ML recipe
2. Average treatment effect on the treated
3. Instrumental variables
4. Standard errors with influence functions
5. More Double ML

## The Double ML recipe

---

# Some definitions

Let

- $O$  be a collection of **observable variables**, e.g.  $O = (W, X, Y)$
- $\theta$  be the **target parameter**
- $\eta$  be the collection of **nuisance parameters**, e.g.  $\eta = (m(X), e(X))$

Double ML uses **score functions**  $\psi(O; \tilde{\theta}, \tilde{\eta})$  that satisfy

1.  $\overbrace{\mathbb{E}[\psi(O; \theta, \eta)]}^{\text{moment condition}} = 0$ , i.e. with expectation zero if evaluated at true parameters
2.  $\partial_r \mathbb{E}[\psi(O; \theta, \eta + r(\tilde{\eta} - \eta))]|_{r=0}$ , i.e. **Neyman-orthogonality**

## Examples

Recall that the moment condition of the residual-on-residual regression with  $m(X) \equiv \mathbb{E}[Y | X]$  and  $e(X) \equiv \mathbb{E}[W | X]$  reads:

$$\mathbb{E} [(Y - m(X) - \tau(W - e(X)))(W - e(X))] = 0$$

$$\Rightarrow O = (W, X, Y), \theta = \tau, \eta = (m(X), e(X))$$

AIPW for ATE moment condition with  $m(w, X) \equiv \mathbb{E}[Y | W = w, X]$ :

$$\mathbb{E} \left[ m(1, X) - m(0, X) + \frac{W(Y - m(1, X))}{e(X)} - \frac{(1 - W)(Y - m(0, X))}{1 - e(X)} - \tau_{ATE} \right] = 0$$

$$\Rightarrow O = (W, X, Y), \theta = \tau_{ATE}, \eta = (m(1, X), m(0, X), e(X))$$

# Linear score functions

We will focus on linear score functions that can be represented as

$$\psi(O; \tilde{\theta}, \tilde{\eta}) = \tilde{\theta} \psi_a(O; \tilde{\eta}) + \psi_b(O; \tilde{\eta})$$

such that the moment condition can be written as

$$\mathbb{E}[\psi(O; \theta, \eta)] = \theta \mathbb{E}[\psi_a(O; \eta)] + \mathbb{E}[\psi_b(O; \eta)] = 0$$

and the solution is

$$\theta = -\frac{\mathbb{E}[\psi_b(O; \eta)]}{\mathbb{E}[\psi_a(O; \eta)]}$$

## Example residual-on-residual regression

Moment condition:

$$\begin{aligned}\mathbb{E} [(Y - m(X) - \tau(W - e(X)))(W - e(X))] &= 0 \\ \mathbb{E} [(Y - m(X))(W - e(X)) - \tau(W - e(X))(W - e(X))] &= 0 \\ \tau \underbrace{\mathbb{E}[(-1)(W - e(X))^2]}_{\psi_a} + \underbrace{\mathbb{E}[(Y - m(X))(W - e(X))]}_{\psi_b} &= 0 \\ \Rightarrow \tau &= -\frac{\mathbb{E}[\psi_b(O; \eta)]}{\mathbb{E}[\psi_a(O; \eta)]} = \frac{\mathbb{E} [(Y - m(X))(W - e(X))]}{\mathbb{E} [(W - e(X))^2]}\end{aligned}$$



## Example AIPW

AIPW for ATE moment condition:

$$\begin{aligned} \mathbb{E} \left[ m(1, X) - m(0, X) + \frac{W(Y - m(1, X))}{e(X)} - \frac{(1 - W)(Y - m(0, X))}{1 - e(X)} - \tau_{ATE} \right] &= 0 \\ \underbrace{\tau_{ATE}(-1)}_{\psi_a} + \mathbb{E} \left[ \underbrace{m(1, X) - m(0, X) + \frac{W(Y - m(1, X))}{e(X)} - \frac{(1 - W)(Y - m(0, X))}{1 - e(X)}}_{\psi_b} \right] &= 0 \\ \Rightarrow \tau_{ATE} = -\frac{\mathbb{E}[\psi_b(O; \eta)]}{\mathbb{E}[\psi_a(O; \eta)]} = \mathbb{E} \left[ m(1, X) - m(0, X) + \frac{W(Y - m(1, X))}{e(X)} - \frac{(1 - W)(Y - m(0, X))}{1 - e(X)} \right] \end{aligned}$$

This is a very complicated way to say that we take the expectation of the pseudo-outcome we called  $\tilde{Y}_{ATE}$  last week, but it illustrates the recipe

## Double ML recipe

1. Find **Neyman-orthogonal score** for your target parameter (they can be constructed, see Sec. 2 of [Chernozhukov et al., 2018](#))
2. **Predict nuisance parameters**  $\hat{\eta}$  with cross-fitted high-quality ML
3. **Solve empirical moment condition** to estimate the target parameter

$$\hat{\theta} = -\frac{\sum_i \psi_b(O_i; \hat{\eta}_i)}{\sum_i \psi_a(O_i; \hat{\eta}_i)}$$

4. **Calculate standard error**

$$\hat{\sigma}^2 = \frac{N^{-1} \sum_i \psi(O_i; \hat{\theta}, \hat{\eta}_i)^2}{[N^{-1} \sum_i \psi_a(O_i; \hat{\eta}_i)]^2} \Rightarrow se(\hat{\theta}) = \sqrt{\frac{\hat{\sigma}^2}{N}}$$



Don't panic, we will unpack this later, but first some use cases

Average treatment effect on the  
treated

---

# Average treatment effect on the treated

Recall that we defined also the

AVERAGE TREATMENT EFFECT ON THE TREATED (ATT):  $\tau_{ATT} \equiv \mathbb{E}[Y(1) - Y(0) \mid W = 1]$

- What is the expected treatment effect in the subpopulation actually receiving the treatment?

Why is this an interesting parameter?

It helps us to evaluate the quality of treatment assignment if we compare it to the *ATE* (assuming higher outcomes are better):

- $ATT > ATE$ : Treatment assignment better than random
- $ATT = ATE$ : Treatment assignment as good as random
- $ATT < ATE$ : Treatment assignment worse than random

# ATT AIPW moment condition

This is the Neyman-orthogonal moment condition for ATT

$$\mathbb{E} \left[ \underbrace{\frac{W}{e}(Y - m(0, X))}_{\text{regression adjustment}} - \underbrace{\frac{(1 - W)e(X)}{e(1 - e(X))}(Y - m(0, X))}_{\text{IPW weighted residual}} - \tau_{ATT} \frac{W}{e} \right] = 0 \quad (1)$$

where  $e = \mathbb{E}[W]$

$\Rightarrow O = (W, X, Y), \theta = \tau_{ATT}, \eta = (m(0, X), e(X))^1, \psi_a = (-1)\frac{W}{e}, \psi_b = \tilde{Y}_{ATT}$

---

<sup>1</sup> $e$  is a constant and not a nuisance parameter

## Instrumental variables

---

## The framework: graph unconditional

A very popular research design (at least in economics) for identifying causal effects is to assume access to an **instrumental variable**

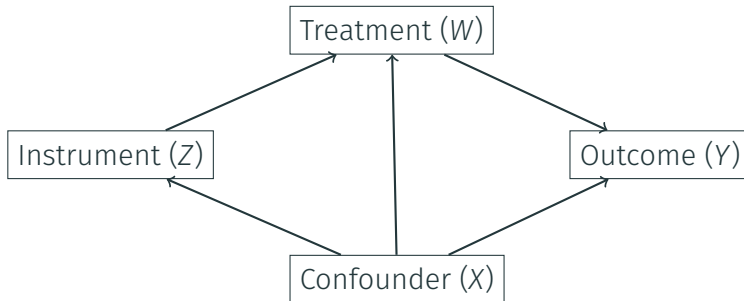


A valid instrument affects the outcome only through the treatment and we can **use this exogenous treatment variation** to identify its effect

Nice refresher

## The framework: graph conditional

If  $Z$  is not 🎲, we still need to adjust for confounders



Current workhorse for estimation is 2SLS, but same model selection issues like for OLS



# What if partially linear model identification is not credible

## Problem:

If identifying assumption 1 is not fulfilled, we expect

$$\mathbb{E}[Y \mid W, X] = \tau W + g(X) + \underbrace{\mathbb{E}[U_{Y_W} \mid W, X]}_{\neq 0}$$

Identification ❌

⇒ Plain partially linear estimator would be biased 😱

If we are lucky, we can leverage the exogenous variation induced by instrumental variable  $Z$

## Partially linear IV model

Under modelling assumption 2 that the potential outcome is partially linear and

### Identifying Assumption 3 (conditionally exogenous instrument)

The potential outcomes specific error term is conditionally independent of  $Z$ :

$$Y(w) = \tau w + g(X) + U_{Y_w}; \quad \mathbb{E}[U_{Y_w} | Z, X] = 0; \quad \forall w \in \mathcal{W}$$

we can identify   $\tau$  as follows, where  $h(X) \equiv \mathbb{E}[Z | X]$

$$\tau = \frac{\mathbb{E}[(Z - h(X))(Y - m(X))]}{\mathbb{E}[(Z - h(X))(W - e(X))]} = \frac{\text{Cov}[Z - h(X), Y - m(X)]}{\text{Cov}[Z - h(X), W - e(X)]}$$

$\Rightarrow$  using residuals of predicted instrument as instrument for the residual-on-residual regression

# Identification in partially linear IV model

Recall that  $Y = \tau W + g(X) + U_{YW}$  and  $m(X) = \tau e(X) + g(X)$  and suppress dependencies of nuisance parameters on  $X$

$$\begin{aligned}\tau &= \frac{\text{Cov}[Z - h, Y - m]}{\text{Cov}[Z - h, W - e]} = \frac{\text{Cov}[Z - h, \tau W + g + U_{YW} - \tau e - g]}{\text{Cov}[Z - h, W - e]} \\ &= \frac{\text{Cov}[Z - h, \tau(W - e) + U_{YW}]}{\text{Cov}[Z - h, W - e]} \\ &= \tau \frac{\text{Cov}[Z - h, W - e]}{\text{Cov}[Z - h, W - e]} + \frac{\text{Cov}[Z - h, U_{YW}]}{\text{Cov}[Z - h, W - e]} \\ &= \tau + \frac{\text{Cov}[Z, U_{YW}]}{\text{Cov}[Z - h, W - e]} - \frac{\text{Cov}[h, U_{YW}]}{\text{Cov}[Z - h, W - e]} \\ &\stackrel{\text{IA3}}{=} \tau\end{aligned}$$

$$\text{b/c } \text{Cov}[Z, U_{YW}] = \mathbb{E}[ZU_{YW}] - \underbrace{\mathbb{E}[Z] \mathbb{E}[U_{YW}]}_{=0} \stackrel{\text{LIE}}{=} \mathbb{E}[\mathbb{E}[ZU_{YW} \mid Z, X]] = \mathbb{E}[Z \underbrace{\mathbb{E}[U_{YW} \mid Z, X]}_{\stackrel{\text{IA3}}{=}0}] = 0, \text{ same for}$$

$$\text{Cov}[h(X), U_{YW}] = 0$$

Note that replacing  $Z$  by  $W$  in all equations and IA3 recovers identification of the partially linear model w/o IV

## Partially linear IV moment condition

Robinson/partialling out style **moment condition with Neyman-orthogonal score**:<sup>2</sup>

$$\begin{aligned} \mathbb{E} \left[ \begin{pmatrix} \overbrace{Y - m(X)}^{\text{outcome residual}} & -\tau & \overbrace{(W - e(X))}^{\text{treatment residual}} & \overbrace{(Z - h(X))}^{\text{instrument residual}} \end{pmatrix} \right] = 0 \\ \mathbb{E} [(Y - m(X))(Z - h(X)) - \tau(W - e(X))(Z - h(X))] = 0 \\ \tau \underbrace{\mathbb{E}[(-1)(W - e(X))(Z - h(X))]}_{\psi_a} + \underbrace{\mathbb{E}[(Y - m(X))(Z - h(X))]}_{\psi_b} = 0 \end{aligned}$$

$$\Rightarrow O = (W, X, Y, Z), \theta = \tau, \eta = (m(X), e(X), h(X))$$

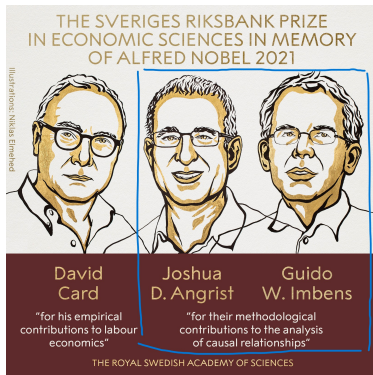
The standard recipe applies (implemented in **DoubleML**)

---

<sup>2</sup>You find alternative Neyman-orthogonal scores in Section 4.2 of **Chernozhukov et al. (2018)**.

# Nobel price alert

What if effects are not homogeneous?



Angrist explanation

Imbens explanation

## Local Average Treatment Effects (1/3)

In the special case of a randomized binary instrument  $Z$  and a binary treatment  $W(Z)$ , Imbens & Angrist (1994) show that the Wald estimator

$$\frac{\overbrace{\mathbb{E}[Y \mid Z = 1] - \mathbb{E}[Y \mid Z = 0]}^{\text{reduced form / intention to treat}}}{\underbrace{\mathbb{E}[W \mid Z = 1] - \mathbb{E}[W \mid Z = 0]}_{\text{first stage / complier share}}} = E[Y(1) - Y(0) \mid \underbrace{W(1) - W(0) = 1}_{\text{complier}}] \equiv \tau_{LATE}$$

identifies a LOCAL AVERAGE TREATMENT EFFECT (LATE) under monotonicity that nobody is moved out of the treatment by the instrument (no defiers)

The LATE describes the effect of the compliers who change their treatment status due to the instrument (may or may not be an interesting target parameter)

## Local Average Treatment Effects (2/3)

What if  $Z$  is not randomly assigned, but we assume we observe all confounders?  
(see [Frölich \(2007\)](#) for identification details)

Double ML uses the following moment condition with a **Neyman-orthogonal score**

$$\mathbb{E} \left[ \overbrace{m_z(1, X) - m_z(0, X) + \frac{Z(Y - m_z(1, X))}{h(X)} - \frac{(1 - Z)(Y - m_z(0, X))}{1 - h(X)}}^{\psi_b} \right] \quad (2)$$
$$+ \mathbb{E} \left[ \underbrace{(-1) \left[ e(1, X) - e(0, X) + \frac{Z(W - e(1, X))}{h(X)} - \frac{(1 - Z)(W - e(0, X))}{1 - h(X)} \right]}_{\psi_a} \right] \times \tau_{LATE} = 0$$

where  $m_z(z, X) = \mathbb{E}[Y \mid Z = z, X]$ ,  $h(X) = \mathbb{P}(Z = 1 \mid X)$  is the probability to be “instrumented” and  $e(Z, X) = \mathbb{P}(W = 1 \mid Z, X)$  the probability to be treated given the instrument

## Local Average Treatment Effects (3/3)

Equation (2) looks terrifying but leads to a **familiar structure**

$$\tau_{LATE} = \frac{\overbrace{\mathbb{E} \left[ m_z(1, X) - m_z(0, X) + \frac{Z(Y - m_z(1, X))}{h(X)} - \frac{(1 - Z)(Y - m_z(0, X))}{1 - h(X)} \right]}^{\text{reduced form / intention to treat}}}{\underbrace{\mathbb{E} \left[ e(1, X) - e(0, X) + \frac{Z(W - e(1, X))}{h(X)} - \frac{(1 - Z)(W - e(0, X))}{1 - h(X)} \right]}_{\text{first stage / complier share}}} \quad (3)$$

It **generalizes the Wald estimator** to the case with confounders

It just **divides the ATE of the instrument on the outcome** (reduced form) by the **ATE of the instrument on the treatment** (first stage)



## Standard errors with influence functions

---

# How to do statistical inference for Double ML

I told you to **estimate the standard error** like this

$$\hat{\sigma}^2 = \frac{N^{-1} \sum_i \psi(O_i; \hat{\theta}, \hat{\eta}_i)^2}{[N^{-1} \sum_i \psi_a(O_i; \hat{\eta}_i)]^2} \Rightarrow se(\hat{\theta}) = \sqrt{\frac{\hat{\sigma}^2}{N}}$$

But why? 🤔

For better understanding, we need to introduce the **concept of influence functions**

Influence functions are powerful tools beyond Double ML, but we will focus on its use for our special case with linear scores and do not go into the technical details

For more general introductions see Kahn (2022), Jann (2019, 2020) or this [tweeterial](#) for a short intro

# Influence function

How is it defined?

$$\Psi(O; \theta, \eta) = -\mathbb{E}\left[\frac{\delta\psi}{\delta\theta}\right]^{-1} \psi(O; \theta, \eta) = -\mathbb{E}[\psi_a(O; \eta)]^{-1} \psi(O; \theta, \eta)$$

⇒ It is a **scaled version of the score** evaluated at the true parameter values

Important features:

- $\mathbb{E}[\Psi(O; \theta, \eta)] = \mathbb{E}[-\mathbb{E}[\psi_a(O; \eta)]^{-1} \psi(O; \theta, \eta)] = -\mathbb{E}[\psi_a(O; \eta)]^{-1} \underbrace{\mathbb{E}[\psi(O; \theta, \eta)]}_{=0} = 0$
- $\Psi(O_i; \theta, \eta_i)/N$  **approximates the influence of observation  $i$**  on the estimate of the target parameter

## Influence function for inference (1/2)

What makes it so valuable for statistical inference?

$$\sqrt{N}(\hat{\theta} - \theta) = \frac{1}{\sqrt{N}} \sum_i \Psi(O_i; \theta, \eta_i) + o_p(1) \xrightarrow{d} N(0, \underbrace{\text{Var}[\Psi(O; \theta, \eta)]}_{\sigma^2})$$

⇒ The estimator distribution and the influence function are closely linked

Note that (suppressing the arguments)

$$\sigma^2 = \text{Var}[\Psi] = \mathbb{E}[\Psi^2] - \underbrace{\mathbb{E}[\Psi]^2}_{=0} = \mathbb{E}[\Psi^2] = \mathbb{E}[(-\mathbb{E}[\psi_a]^{-1}\psi)^2] = \mathbb{E}[\psi_a]^{-2} \mathbb{E}[\psi^2] = \frac{\mathbb{E}[\psi^2]}{\mathbb{E}[\psi_a]^2}$$

## Influence function for inference (2/2)

The **sample equivalent** is therefore

$$\hat{\sigma}^2 = \frac{N^{-1} \sum_i \psi(O_i; \hat{\theta}, \hat{\eta}_i)^2}{[N^{-1} \sum_i \psi_a(O_i; \hat{\eta}_i)]^2}$$

and the standard error is calculated as  $se(\hat{\theta}) = \sqrt{\frac{\hat{\sigma}^2}{N}}$

$se(\hat{\theta})$  can then be used to calculate **t-values, confidence intervals** etc.

## Example ATE

Start with the ATE score where we denote  $m(w, X) = m_w$  and  $e(X) = e$

$$\psi_{ATE} = \underbrace{m_1 - m_0 + \frac{(D - e)(Y - m_W)}{e(1 - e)}}_{\psi_b} + \underbrace{(-1)}_{\psi_a} \tau_{ATE}$$

$$\Rightarrow \frac{\delta \psi_{ATE}}{\delta \tau_{ATE}} = \psi_a = -1$$

$$\Rightarrow \Psi_{ATE} = -\mathbb{E}[(-1)]^{-1}[\tilde{Y}_{ATE} - \tau_{ATE}] = \tilde{Y}_{ATE} - \tau_{ATE}$$

In the special case where  $\mathbb{E}[\psi_a] = -1$ , score and influence function coincide

## Influence function for inference (bonus)

A cool/convenient feature of influence functions is that they obey the **chain rule**

Imagine that your target parameter is a function of  $k$  other parameters, i.e.

$$\theta = f(\theta_1, \dots, \theta_K)$$

Then the influence function of  $\theta$  is

$$\psi_{\theta} = \sum_{k=1}^K \frac{\delta f}{\delta \theta_k} \psi_{\theta_k}$$

$\Rightarrow$  we can use existing influence functions to create new ones

## Example $ATT - ATE$

We may want to test whether the  $ATT = ATE$

For this purpose, we create the new parameter  $\Delta(\tau_{ATT}, \tau_{ATE}) = \tau_{ATT} - \tau_{ATE}$

The new influence function is

$$\begin{aligned}\psi_{\Delta} &= \overbrace{\frac{\delta \Delta}{\delta \tau_{ATT}}}^{=1} \psi_{\tau_{ATT}} + \overbrace{\frac{\delta \Delta}{\delta \tau_{ATE}}}^{=-1} \psi_{\tau_{ATE}} \\ &= \psi_{\tau_{ATT}} - \psi_{\tau_{ATE}} \\ &= \tilde{Y}_{ATT} - \tau_{ATT}W/e - \tilde{Y}_{ATE} + \tau_{ATE}\end{aligned}$$

and can be used to get  $se(\hat{\Delta})$



## Two ways to get the *LATE* influence function: (1) direct way

The direct way starts with the score implied by (2) where we denote  $m(z, X) = m_z$ ,  $e(z, X) = e_z$ ,  $h(X) = h$ , and rewrite the weighted residuals for compactness

$$\psi_{LATE} = \underbrace{m_1 - m_0 + \frac{\tilde{Y}_{Z \rightarrow Y}(Z - h)(Y - m_Z)}{h(1 - h)}}_{\psi_b} - \underbrace{\left[ e_1 - e_0 + \frac{\tilde{Y}_{Z \rightarrow W}(Z - h)(Y - e_Z)}{h(1 - h)} \right]}_{\psi_a} \times \tau_{LATE}$$

$$\Rightarrow \frac{\delta \psi_{LATE}}{\delta \tau_{LATE}} = \psi_a = -\tilde{Y}_{Z \rightarrow W}$$

$$\begin{aligned} \Rightarrow \Psi_{LATE} &= -\mathbb{E}[-\tilde{Y}_{Z \rightarrow W}]^{-1} [\tilde{Y}_{Z \rightarrow Y} - \tilde{Y}_{Z \rightarrow W} \times \tau_{LATE}] \\ &= \mathbb{E}[\tilde{Y}_{Z \rightarrow W}]^{-1} [\tilde{Y}_{Z \rightarrow Y} - \tilde{Y}_{Z \rightarrow W} \times \tau_{LATE}] \end{aligned}$$

## Two ways to get the $LATE$ influence function: (2) chain rule

Denote by  $\tau_{Z \rightarrow W}$  and  $\tau_{Z \rightarrow Y}$  the ATEs of  $Z$  on  $W$  and  $Y$ , respectively

We start by noting that

$$\tau_{LATE}(\tau_{Z \rightarrow Y}, \tau_{Z \rightarrow W}) = \frac{\tau_{Z \rightarrow Y}}{\tau_{Z \rightarrow W}}$$

The chain rule tells us that

$$\begin{aligned}\psi_{LATE} &= \frac{\delta \tau_{LATE}}{\delta \tau_{Z \rightarrow Y}} \psi_{\tau_{Z \rightarrow Y}} + \frac{\delta \tau_{LATE}}{\delta \tau_{Z \rightarrow W}} \psi_{\tau_{Z \rightarrow W}} = \frac{1}{\tau_{Z \rightarrow W}} \psi_{\tau_{Z \rightarrow Y}} - \frac{\tau_{Z \rightarrow Y}}{\tau_{Z \rightarrow W}^2} \psi_{\tau_{Z \rightarrow W}} \\&= \frac{1}{\tau_{Z \rightarrow W}} (\tilde{Y}_{Z \rightarrow Y} - \tau_{Z \rightarrow Y}) - \frac{\tau_{Z \rightarrow Y}}{\tau_{Z \rightarrow W}^2} (\tilde{Y}_{Z \rightarrow W} - \tau_{Z \rightarrow W}) \\&= \frac{1}{\tau_{Z \rightarrow W}} \left[ \tilde{Y}_{Z \rightarrow Y} - \cancel{\tau_{Z \rightarrow Y}} - \underbrace{\frac{\tau_{Z \rightarrow Y}}{\tau_{Z \rightarrow W}} \tilde{Y}_{Z \rightarrow W}}_{\tau_{LATE}} + \cancel{\frac{\tau_{Z \rightarrow Y}}{\tau_{Z \rightarrow W}} \tau_{Z \rightarrow W}} \right] \\&= \mathbb{E}[\tilde{Y}_{Z \rightarrow W}]^{-1} [\tilde{Y}_{Z \rightarrow Y} - \tilde{Y}_{Z \rightarrow W} \times \tau_{LATE}]\end{aligned}$$

$$\text{b/c } \tau_{Z \rightarrow W} = \mathbb{E}[\tilde{Y}_{Z \rightarrow W}]$$

Application notebook: Double ML as generic recipe

## More Double ML

---

# The generic concept is applied to many other scenarios

**2x2 Difference-in-differences:** Chang (2020) and Zimmert (2020)

But **not so much for the more general panel** case (this literature is currently productively defining what parameters canonical estimators imply #LATEreloaded)

I am convinced Double ML estimators in this direction are currently work in progress (or I missed them)

**Mediation:** Farbmacher et al. (2022)

**Dynamic treatments:** Bodory et al. (2022)

**Quantile treatment effects:** Belloni et al. (2017) and Kallus, Mao & Uehara (2019)  
and many more (to come)

Important for you is that they build on the **same/similar principles**

*Ceterum censeo* a fancy method alone is not a credible  
identification strategy  
⇒ separate identification and estimation