

# Treatment Effect Estimators as Weighted Outcomes

Michael C. Knaus\*

December 13, 2024

## Abstract

Estimators that weight observed outcomes to form effect estimates have a long tradition. Their outcome weights are widely used in established procedures, such as checking covariate balance, characterizing target populations, or detecting and managing extreme weights. This paper introduces a general framework for deriving such outcome weights. It establishes when and how numerical equivalence between an original estimator representation as moment condition and a unique weighted representation can be obtained. The framework is applied to derive novel outcome weights for the six seminal instances of double machine learning and generalized random forests, while recovering existing results for other estimators as special cases. The analysis highlights that implementation choices determine (i) the availability of outcome weights and (ii) their properties. Notably, standard implementations of partially linear regression-based estimators, like causal forests, employ outcome weights that do not sum to (minus) one in the (un)treated group, not fulfilling a property often considered desirable.

**Keywords:** Augmented inverse probability weighting, causal forest, causal machine learning, covariate balancing, double machine learning, generalized random forest, implied weights, partially linear regression

---

\*University of Tübingen, Mohlstraße 36, 72074 Tübingen, Germany. Michael C. Knaus is also affiliated with IZA, Bonn, [michael.knaus@uni-tuebingen.de](mailto:michael.knaus@uni-tuebingen.de). I thank participants of seminars at Aarhus, Bologna, CREST, Düsseldorf, LISER, Mannheim, and St. Gallen, of EuroCIM 2024 and VfS 2024, and in particular Oliver Dukes, Noah Greifer, Phillip Heiler, Michael Lechner, Dor Leventer, David Preinerstorfer, Christoph Rothe, Tymon Słoczyński, Anthony Strittmatter, Linbo Wang and Jose Zubizarreta for valuable comments improving the paper. I thank Henri Pfleiderer for assisting with the implementation of the R package. The usual disclaimer applies.

# 1 Introduction

Estimating the effect of treatment  $D_i$  on outcome  $Y_i$  is a common goal in causal inference. A variety of estimators is available for estimating different target parameters, after arguing for their identification within a particular research design (see, e.g. reviews by [Imbens & Wooldridge, 2009](#); [Athey & Imbens, 2017](#); [Abadie & Cattaneo, 2018](#); [Imbens, 2024](#)). Many of these estimators are a “white box” in the sense that they document how the sample is processed to obtain an effect estimate. Parametric regressions come with familiar coefficient outputs and other popular estimators have a representation as linear combination of observed outcomes:

$$\hat{\tau} = \sum_{i=1}^N \omega_i Y_i = \underbrace{\boldsymbol{\omega}'}_{1 \times N} \underbrace{\mathbf{Y}}_{N \times 1} \quad (1)$$

where  $\omega_i$  represents the weight assigned to the outcome of observation  $i$  in estimating  $\hat{\tau}$ . Structure (1) is most prominent in the literature on propensity score matching/weighting (e.g. [Imbens & Rubin, 2015](#)), balancing estimators (e.g. [Ben-Michael, Feller, Hirshberg, & Zubizarreta, 2021](#)), and synthetic controls (e.g. [Abadie, 2021](#)). Furthermore, it is discussed for estimators of the local average treatment effect (e.g. [Imbens & Rubin, 1997](#); [Abadie, 2003](#); [Śłoczyński, Uysal, & Wooldridge, 2024](#)) as well as for linear regression (e.g. [Imbens, 2015](#); [Chattopadhyay & Zubizarreta, 2023](#)).

Outcome weights  $\omega_i$  have established use cases, such as: (i) *covariate balancing checks* assessing internal validity in experimental and observational studies (e.g. [Rosenbaum & Rubin, 1984, 1985](#)), (ii) *target population characterization* investigating external validity in IV settings (e.g. [Abadie, 2003](#)), or when using OLS ([Chattopadhyay & Zubizarreta, 2023](#)), (iii) *extrapolation diagnostics* for estimators that could use negative weights ([Chattopadhyay & Zubizarreta, 2023](#)), (iv) *finite sample estimator stabilization* by normalizing weights (e.g. [Hájek, 1971](#)), or by trimming extreme weights (e.g. [Lechner & Strittmatter, 2019](#)), (v) *variance estimation* (e.g. [Imbens & Rubin, 2015](#), Ch. 19).

In contrast, recent estimators integrating supervised machine learning into the estimation process (see [Chernozhukov, Hansen, Kallus, Spindler, & Syrgkanis, 2024](#), for a textbook) can be considered as “grey box”. Their multi-step algorithms are transparent

and their theoretical properties are well-understood. However, neither coefficients nor outcome weights are currently available to interrogate how these steps jointly process a concrete sample within a concrete implementation. A key take-away of the analysis below is that at least outcome weights can be available for such multi-step estimators.

This paper introduces a simple but general framework to derive and analyze outcome weights of form (1). We establish conditions for a numerical equivalence between an original estimator representation as moment condition and a unique weighted representation. The framework is applied to derive novel outcome weights for the six seminal instances of double machine learning (Chernozhukov et al., 2018) and generalized random forest (Athey, Tibshirani, & Wager, 2019). Knowing the closed-form of the outcome weights has the immediate practical benefit that they can be plugged into established routines for classic weighting estimators. For example, covariate balancing can now be assessed for conditional average treatment effects estimated by causal forest. A second benefit is that the framework naturally allows to investigate basic properties of the weights. In particular, it highlights that implementation decisions control whether outcome weights of the treated sum up to one and of the untreated to minus one. Such weights are often considered as intuitive, reasonable and desirable in the literature (e.g. Imbens & Rubin, 2015; Słoczyński et al., 2024). However, the new framework reveals that estimators building on partially linear regression do not satisfy this property in standard implementations.

The paper makes several contributions: (i) It introduces the first general framework to derive outcome weights; (ii) Its application to causal machine learning estimators yields novel outcome weights and provides a blueprint for applying the framework to other estimators; (iii) It illustrates how the new closed-form expressions enable established diagnostic tools from the weighting literature to be integrated off-the-shelf into causal machine learning applications; (iv) The theoretical results about conditions ensuring desirable estimator properties inform implementation decisions and complement the high-level conditions provided in asymptotic analyses; (v) The paper provides an additional piece in the continuing effort to blur the line between outcome weighting and outcome regression methods. Bruns-Smith, Dukes, Feller, and Ogburn (2023) show how weighting estimators can be expressed as regression estimators. This paper goes in the opposite

direction by showing how estimators involving flexible outcome regression can be expressed as weighting estimators; (vi) The accompanying R package `OutcomeWeights` computes the weights presented in the paper for general use (Knaus, 2024). The presented applications rely on this package and can be replicated in a supplementary [Docker image](#).

## 1.1 Related literature

Outcome weights in form of (1) are leveraged as common structure of difference, weighting, subclassification, and matching estimators (Smith & Todd, 2005; Huber, Lechner, & Wunsch, 2013; Imbens & Rubin, 2015, Ch. 19.4.). Similarly, the outcome weights are derived and used for ordinary and weighted least squares based estimators (Kline, 2011; Imbens, 2015; Jakiela, 2021; Chattopadhyay & Zubizarreta, 2023; Hazlett & Shinkre, 2024), two-stage least squares (TSLS) (Chattopadhyay & Zubizarreta, 2021), and augmented inverse probability weighting implemented with (post-selection) OLS outcome regression (Knaus, 2021; Chattopadhyay & Zubizarreta, 2023). This paper shows that a broader class of estimators can have this structure with a particular focus on those incorporating flexible outcome regression, while covering the results in the literature as special cases.

Other types of weights are prominent in the causal inference literature but distinct from the *outcome weights* pursued in this paper. First, *balancing weights* are the result of a tailored optimization problem to achieve covariate balancing of some prespecified form (Graham, Pinto, & Egel, 2012; Hainmueller, 2012; Imai & Ratkovic, 2014; Zubizarreta, 2015; Zhao, 2019; Kallus, 2020; Armstrong & Kolesár, 2021; Heiler, 2022). Thus, balancing weights are an explicit part of such balancing estimators. While balancing weights are a special case of outcome weights, this paper focuses on estimators where the outcome weights play no explicit role but are implicit in the common characterization of the estimation procedure. Chattopadhyay and Zubizarreta (2023) call such weights “implied weights” in the context of OLS. Second, *effect weights* are central to understanding the estimand targeted by a given estimator. The pursued structures in this literature are variations of  $\mathbb{E}[w(X_i)\tau(X_i)]$  where  $w(X_i)$  is the weight an estimator assigns to the conditional treatment effect  $\tau(X_i)$  in expectation. Effect weights are derived under different identifying and functional form assumptions for OLS (Angrist, 1998; Angrist & Krueger, 1999; Humphreys,

2009; Aronow & Samii, 2016; Goldsmith-Pinkham, Hull, & Kolesár, 2021; Słoczyński, 2022), TSLS (Imbens & Angrist, 1994; Angrist & Imbens, 1995; Heckman & Vytlacil, 2005; Słoczyński, 2020; Blandhol, Bonney, Mogstad, & Torgovitsky, 2022), two-way fixed-effects (see de Chaisemartin & D’Haultfœuille, 2023; Roth, Sant’Anna, Bilinski, & Poe, 2023, for overviews), regression discontinuity estimators (Lee & Lemieux, 2010), and panel estimators (Chernozhukov, Fernández-Val, Hahn, & Newey, 2013). The main difference between outcome weights and effect weights is that the former apply to observed outcomes and numerically reproduce the estimate without further assumptions, while the latter weight inherently unobservable effects and usually reproduce the estimate only in expectation. Both types of weights have their established use cases and are therefore complementary.

A small but growing body of literature provides nuance to the common notion that it is preferable for the outcome weights to sum to (minus) one within treatment groups. Doudchenko and Imbens (2016) and Breitung, Bolwin, and Töns (2024) challenge this view for synthetic control estimators, and Khan and Ugander (2023) for average treatment effect estimation. Słoczyński et al. (2024) note that some Abadie’s (2003)  $\kappa$  estimators are intermediate cases between normalized and unnormalized estimators, e.g. with treated weights summing to one but untreated weights not to summing minus one. This paper adds the observation that partially linear regression based estimators usually produce weights that overall sum up to zero but not to (minus) one for (un)treated.

Finally, the paper adds to recent works establishing numeric equivalences between different estimator representations (Bruns-Smith et al., 2023) or estimators (Słoczyński, Uysal, & Wooldridge, 2023; Słoczyński et al., 2024) for conceptual and/or practical insights.

## 2 A general framework to derive outcome weights

### 2.1 Notation

The estimators under consideration require access to data with  $N$  observations indexed by  $i = 1, \dots, N$ . The data includes a binary treatment  $D_i$ , an outcome  $Y_i$ , covariates  $\mathbf{X}_i$ , and an optional binary instrument  $Z_i$ , all collected in  $O_i = (D_i, \mathbf{X}_i', Y_i, Z_i)'$ . The empirical

mean of a variable  $A_i$  is represented as  $\mathbb{E}_N[A_i] = N^{-1} \sum_{i=1}^N A_i$ .

Many results are stated in matrix notation where bold letters describe vectors or matrices of variables.  $\mathbf{I}_k$  denotes the identity matrix of dimension  $k$ ,  $\mathbf{0}_k$  and  $\mathbf{1}_k$  represent column vectors of length  $k$  containing zeros and ones, respectively.

## 2.2 Pseudo-IV estimators

This paper focuses on estimators falling into the class of pseudo-IV estimators (PIVE):

**Definition 1** (*pseudo-IV estimators*)

Define the class of pseudo-IV estimators (PIVE) as estimators solving an empirical moment condition of the form

$$\mathbb{E}_N \left[ (\tilde{Y}_i - \hat{\tau} \tilde{D}_i) \tilde{Z}_i \right] = 0 \quad (2)$$

with

- $\tilde{Y}_i = f_Y(O_i; \hat{\eta}_i^Y)$ : scalar pseudo-outcome
- $\tilde{D}_i = f_D(O_i; \hat{\eta}_i^D)$ : scalar pseudo-treatment
- $\tilde{Z}_i = f_Z(O_i; \hat{\eta}_i^Z)$ : scalar pseudo-instrument

where  $\hat{\eta}_i^Y$ ,  $\hat{\eta}_i^D$  and  $\hat{\eta}_i^Z$  are optional nuisance parameters.

Note that the PIVE representation is neither a unique, nor the most compact representation of an estimator. For example, a representation using a linear score  $\mathbb{E}_N \left[ \psi_i^a - \hat{\tau} \psi_i^b \right] = 0$  with  $\psi_i^a = \tilde{Y}_i \tilde{Z}_i$  and  $\psi_i^b = \tilde{D}_i \tilde{Z}_i$  would be equivalent, or vice versa any estimator with a linear score can be written as PIVE with  $\tilde{Y}_i = \psi_i^a$ ,  $\tilde{D}_i = \psi_i^b$ , and  $\tilde{Z}_i = 1$ . However, the PIVE structure is essential for the goal of this paper. In particular, separating the pseudo-outcome from the pseudo-instrument makes the derivation and analysis of the outcome weights tractable.

*Example (OLS):* We use the canonical OLS estimator as a running example to illustrate the general results throughout the paper. Consider a linear outcome model  $Y_i = \tau D_i + \mathbf{X}_i' \boldsymbol{\beta} + \varepsilon$ . The OLS estimator for  $\tau$  can be expressed as PIVE using the residual-on-residual regression representation of the Frisch-Waugh-Lovell Theorem:

$$\mathbb{E}_N \left[ \underbrace{\{Y_i - \mathbf{X}'_i \hat{\beta}_{Y \sim X}\}}_{=: \tilde{Y}_i^{ols}} - \hat{\tau}^{ols} \underbrace{[D_i - \mathbf{X}'_i \hat{\beta}_{D \sim X}]}_{=: \tilde{D}_i^{ols}} \underbrace{[D_i - \mathbf{X}'_i \hat{\beta}_{D \sim X}]}_{=: \tilde{Z}_i^{ols}} \right] = 0 \quad (3)$$

where  $\hat{\beta}_{Y \sim X} := (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$  and  $\hat{\beta}_{D \sim X} := (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{D}$  such that the pseudo-outcome is the outcome residual, and both pseudo-treatment and -instrument are the treatment residual.

### 2.3 Outcome weights of pseudo-IV estimators

Solving Equation 2 leads to parameter estimate

$$\hat{\tau} = \frac{\mathbb{E}_N [\tilde{Z}_i \tilde{Y}_i]}{\mathbb{E}_N [\tilde{Z}_i \tilde{D}_i]} = (\tilde{\mathbf{Z}}' \tilde{\mathbf{D}})^{-1} \tilde{\mathbf{Z}}' \tilde{\mathbf{Y}}. \quad (4)$$

Now assume that the pseudo-outcome vector can be obtained by multiplying a unique  $N \times N$  transformation matrix  $\mathbf{T}$  with the outcome vector, i.e.  $\mathbf{T}\mathbf{Y} = \tilde{\mathbf{Y}}$ . Then, Equation 4 can be written in the form of Equation 1

$$\hat{\tau} = \underbrace{(\tilde{\mathbf{Z}}' \tilde{\mathbf{D}})^{-1} \tilde{\mathbf{Z}}' \mathbf{T}}_{\boldsymbol{\omega}'} \mathbf{Y} = \boldsymbol{\omega}' \mathbf{Y} \quad (5)$$

leading to a core result of the paper:

**Proposition 1** (*outcome weights of PIVE*)

*The outcome weights of a PIVE in the sense of Definition 1 have closed-form*

$$\boldsymbol{\omega}' = (\tilde{\mathbf{Z}}' \tilde{\mathbf{D}})^{-1} \tilde{\mathbf{Z}}' \mathbf{T} \quad (6)$$

*if  $\tilde{\mathbf{Z}}' \tilde{\mathbf{D}} \neq 0$  and a unique transformation matrix  $\mathbf{T}$  exists such that  $\mathbf{T}\mathbf{Y} = \tilde{\mathbf{Y}}$ .*

This simple result is constructive because it motivates a two step procedure to derive outcome weights:

1. Express the estimator as PIVE.
2. Find transformation matrix  $\mathbf{T} \Rightarrow \boldsymbol{\omega}$  has closed-form.

These steps are illustrated below for a variety of estimators. However, the procedure is general and could be pursued for any other estimator fitting into the PIVE structure.

*Example (OLS) continued:* The solution of the residual-on-residual regression (3) in form of Equation 6 is

$$\hat{\tau}^{ols} = \overbrace{\left( \underbrace{\hat{\mathbf{E}}'}_{\tilde{\mathbf{Z}}^{ols'}} \underbrace{\hat{\mathbf{E}}}_{\tilde{\mathbf{D}}^{ols}} \right)^{-1} \underbrace{\hat{\mathbf{E}}'}_{\tilde{\mathbf{Z}}^{ols'}} \underbrace{\mathbf{M}_{\mathbf{X}}}_{\mathbf{T}^{ols}} \mathbf{Y}}^{\omega^{ols'}} \quad (7)$$

where we use the projection matrix  $\mathbf{P}_{\mathbf{X}} := \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$  to define the residual maker matrix  $\mathbf{M}_{\mathbf{X}} := \mathbf{I}_N - \mathbf{P}_{\mathbf{X}}$ , and the treatment residual vector  $\hat{\mathbf{E}} := \mathbf{M}_{\mathbf{X}}\mathbf{D}$ . The residual maker matrix is therefore the outcome transformation matrix of OLS and  $\omega^{ols'} = (\hat{\mathbf{E}}'\hat{\mathbf{E}})^{-1}\hat{\mathbf{E}}'\mathbf{M}_{\mathbf{X}}$  is the outcome weights vector.<sup>1</sup>

### 3 Outcome weights of concrete pseudo-IV estimators

This section leverages the new framework to provide the first characterization of outcome weights for six seminal instances within the double machine learning (DML) and generalized random forest (GRF) frameworks (marked with \*), while also recovering existing results for eight other estimators:

- **IF\***: Instrumental forest (Athey et al., 2019)
- PLR-IV\*: Partially linear regression with IV (Chernozhukov et al., 2018)
- TSLS: Two stage least squares
- Wald: Wald estimator (Wald, 1940)
- CF\*: Causal forest (Athey et al., 2019)
- PLR\*: Partially linear regression (Robinson, 1988; Chernozhukov et al., 2018)
- OLS: Ordinary least squares
- DiM: Difference in means
- **AIPW\***: Augmented inverse probability weighting (Robins & Rotnitzky, 1995; Chernozhukov et al., 2018)
- RA: Regression adjustment (e.g. discussed by Imbens, 2004)

---

<sup>1</sup>We deliberately do not use that  $\mathbf{M}_{\mathbf{X}}$  is idempotent for illustration purposes but note that also the identity matrix would be a suitable transformation matrix in (7).



- IPW: Inverse probability weighting (Horvitz & Thompson, 1952)
- **Wald-AIPW\***: **Wald type AIPW** (Tan, 2006; Chernozhukov et al., 2018)
- Wald-RA: Wald type regression adjustment (Tan, 2006)
- Wald-IPW: Wald type inverse probability weighting (Tan, 2006)

Conveniently it suffices to analyze the three estimators in bold letters - IF, AIPW and Wald-AIPW - because their subsequent estimators follow as special cases (see Figure A.1 for a graphical illustration). Each estimator is typically used at an intersection of three research designs (randomized controlled trials, unconfoundedness or instrumental variables), two aggregation levels (average or conditional effects), and three outcome model assumptions (none, partially linear, or linear models). Appendix A.1 summarizes the causal parameters and settings for which each estimator is usually applied. However, the main text ignores definition, identification and interpretation issues concentrating on the mechanics of the estimators.

### 3.1 Nuisance parameters

#### 3.1.1 Definitions

The considered estimators require a variety of nuisance parameters in the form of approximated conditional expectations:

$$\begin{aligned}
\hat{Y}_i &:= \hat{\mathbb{E}}[Y_i | \mathbf{X}_i] & \hat{D}_i &:= \hat{\mathbb{E}}[D_i | \mathbf{X}_i] \\
\hat{Y}_{0,i}^d &:= \hat{\mathbb{E}}[Y_i | D_i = 0, \mathbf{X}_i] & \hat{D}_{0,i}^z &:= \hat{\mathbb{E}}[D_i | Z_i = 0, \mathbf{X}_i] \\
\hat{Y}_{1,i}^d &:= \hat{\mathbb{E}}[Y_i | D_i = 1, \mathbf{X}_i] & \hat{D}_{1,i}^z &:= \hat{\mathbb{E}}[D_i | Z_i = 1, \mathbf{X}_i] \\
\hat{Y}_{0,i}^z &:= \hat{\mathbb{E}}[Y_i | Z_i = 0, \mathbf{X}_i] & \hat{Z}_i &:= \hat{\mathbb{E}}[Z_i | \mathbf{X}_i] \\
\hat{Y}_{1,i}^z &:= \hat{\mathbb{E}}[Y_i | Z_i = 1, \mathbf{X}_i] & &
\end{aligned} \tag{8}$$

Furthermore, define the inverse probability weights of the treated  $\lambda_{1,i}^{ipw} := D_i / \hat{D}_i$  and of the untreated  $\lambda_{0,i}^{ipw} := (1 - D_i) / (1 - \hat{D}_i)$ .<sup>2</sup> Similarly, define the instrument inverse probability weights as  $\lambda_{1,i}^{ipw,z} := Z_i / \hat{Z}_i$  and  $\lambda_{0,i}^{ipw,z} := (1 - Z_i) / (1 - \hat{Z}_i)$ .

<sup>2</sup>We use  $\lambda$  to remind us that these weights are on a different scale than the  $\omega$  weights. Using the corresponding  $\omega_{d,i}^{ipw} := \lambda_{d,i}^{ipw} / N$  definition would unnecessarily complicate notation below.

### 3.1.2 A crucial building block: Smoothers

The literature knows numerous regression methods to estimate the outcome nuisance parameters  $\hat{Y}_i$ ,  $\hat{Y}_{d,i}^d$ , and  $\hat{Y}_{z,i}^z$ . However, the class of smoothers (see e.g. [Hastie & Tibshirani, 1990](#), Ch. 2-3) turns out to be crucial for the purpose of this paper. Smoothers produce outcome predictions by weighting/smoothing observed outcomes

$$\hat{Y}_i = \sum_{j=1}^N s_{i \leftarrow j} Y_j \quad (9)$$

where the *smoother weight*  $s_{i \leftarrow j}$  represents the contribution of unit  $j$ 's outcome to the prediction of unit  $i$ .<sup>3</sup> Define also the *smoother vector* for the outcome prediction of unit  $i$  by  $\mathbf{s}_i = (s_{i \leftarrow 1}, \dots, s_{i \leftarrow N})'$  and the  $N \times N$  *smoother matrix*  $\mathbf{S} = [\mathbf{s}_1 \dots \mathbf{s}_N]'$  such that  $\mathbf{s}'_i \mathbf{Y} = \hat{Y}_i$  and  $\mathbf{S} \mathbf{Y} = \hat{\mathbf{Y}}$ .

The smoother weights in this paper are explicitly allowed to depend on the outcomes (adaptive smoother) and on random components (random smoother), i.e.  $\mathbf{s}_i := \mathbf{s}_i(\mathbf{X}_i; \mathbf{X}, \mathbf{Y}, \epsilon_s)$ .<sup>4</sup> This covers for example (post-selection) OLS, ridge, spline and kernel (ridge) regressions, regression trees, random forests or boosted trees with data-driven hyperparameter tuning (see [Appendix A.2](#) for further discussion). However, the numerical equivalences established below require the mere existence of a smoother matrix:

**Condition 1** (*smoother matrix*)

*There exists a unique smoother matrix creating the outcome nuisance vectors if multiplied with the outcome vector:*

$$(C1a) \quad \mathbf{S} \mathbf{Y} = \hat{\mathbf{Y}}$$

$$(C1b) \quad \mathbf{S}_0^d \mathbf{Y} = \hat{\mathbf{Y}}_0^d \quad \text{and} \quad \mathbf{S}_1^d \mathbf{Y} = \hat{\mathbf{Y}}_1^d$$

$$(C1c) \quad \mathbf{S}_0^z \mathbf{Y} = \hat{\mathbf{Y}}_0^z \quad \text{and} \quad \mathbf{S}_1^z \mathbf{Y} = \hat{\mathbf{Y}}_1^z$$

*Example (OLS) continued:* The projection matrix is arguably the most prominent smoother matrix producing fitted values of an OLS regression as  $\underbrace{\mathbf{P}_X}_{\mathbf{S}^{ols}} \mathbf{Y} = \hat{\mathbf{Y}}^{ols}$ .

<sup>3</sup>The arrow notation is adapted from [Lin and Han \(2022\)](#).

<sup>4</sup>The categorization of smoothers is inspired by [Curth, Jeffares, and van der Schaar \(2024\)](#).

## 3.2 Concrete outcome weights

### 3.2.1 Instrumental forest and its special cases

The instrumental forest (IF) of [Athey et al. \(2019\)](#) runs an  $\mathbf{x}$ -specific weighted partially linear IV regression

$$\mathbb{E}_N \left[ \underbrace{\{Y_i - \hat{Y}_i\}}_{=: \tilde{Y}_i^{if}} - \hat{\tau}^{if}(\mathbf{x}) \underbrace{[D_i - \hat{D}_i]}_{=: \tilde{D}_i^{if}} \underbrace{[Z_i - \hat{Z}_i]}_{=: \tilde{Z}_i^{if}} \alpha_i^{if}(\mathbf{x}) \right] = 0 \quad (10)$$

where the  $\mathbf{x}$ -specific weights  $\alpha^{if}(\mathbf{x})$  are obtained by the tailored splitting criterion described in [Athey et al. \(2019\)](#) and can be extracted via the `get_forest_weights()` function of their `grf` R package ([Tibshirani, Athey, Sverdrup, & Wager, 2024](#)). The solution in the form of Equation 4 is

$$\hat{\tau}^{if}(\mathbf{x}) = (\hat{\mathbf{R}}' \text{diag}(\boldsymbol{\alpha}^{if}(\mathbf{x})) \hat{\mathbf{V}})^{-1} \hat{\mathbf{R}}' \text{diag}(\boldsymbol{\alpha}^{if}(\mathbf{x})) \hat{\mathbf{U}} \quad (11)$$

where  $\hat{\mathbf{R}} = \mathbf{Z} - \hat{\mathbf{Z}}$ ,  $\hat{\mathbf{V}} = \mathbf{D} - \hat{\mathbf{D}}$  and  $\hat{\mathbf{U}} = \mathbf{Y} - \hat{\mathbf{Y}}$  are the instrument, treatment and outcome residual vectors, respectively. The PIVE structure is therefore established. The next step is to understand whether the pseudo-outcome can be obtained using a transformation matrix. This is only possible if a smoother is applied to obtain the outcome predictions, i.e. Condition 1a holds such that  $\hat{\mathbf{U}} = \mathbf{Y} - \mathbf{S}\mathbf{Y} = (\mathbf{I}_N - \mathbf{S})\mathbf{Y}$  and

$$\hat{\tau}^{if}(\mathbf{x}) = \underbrace{(\hat{\mathbf{R}}' \text{diag}(\boldsymbol{\alpha}^{if}(\mathbf{x})) \hat{\mathbf{V}})^{-1} \hat{\mathbf{R}}' \text{diag}(\boldsymbol{\alpha}^{if}(\mathbf{x}))}_{\boldsymbol{\omega}^{if'}} \overbrace{(\mathbf{I}_N - \mathbf{S})}^{\mathbf{T}^{if}} \mathbf{Y}. \quad (12)$$

The transformation matrix of IF can therefore be considered as a generalized residual maker matrix. Equation 12 contains then the first concrete case of Proposition 1:

**Corollary 1** (*outcome weights of instrumental forest*)

*Under Condition 1a such that the outcome predictions can be written as  $\mathbf{S}\mathbf{Y} = \hat{\mathbf{Y}}$ , the outcome weights of instrumental forests take the form*

$$\hat{\omega}^{if}(\mathbf{x})' = (\hat{\mathbf{R}}' \text{diag}(\boldsymbol{\alpha}^{if}(\mathbf{x})) \hat{\mathbf{V}})^{-1} \hat{\mathbf{R}}' \text{diag}(\boldsymbol{\alpha}^{if}(\mathbf{x})) (\mathbf{I}_N - \mathbf{S}). \quad (13)$$

Table 1 compactly shows how seven other estimators (in light gray) follow as special cases of IF. Starting from the dark gray row, we can follow an upward path to the Wald estimator or a downward path to DiM. The white rows between the gray rows document the modifications needed to recover the next estimator. For example, moving from IF to CF uses treatment residuals instead of instrument residuals and the CF specific weights  $\alpha^{cf}$  in the pseudo-instrument, while pseudo-treatment and transformation matrix remain unchanged. Similarly setting the weights to one recovers PLR from CF and PLR-IV from IF. Continuing the paths up- and downwards replaces the generic predictions with linear projections to recover TSLS and OLS, respectively. Finally, using the projection matrix of a constant recovers Wald estimator and DiM.

Table 1: IF pseudo-variables and transformation matrices

	$\tilde{\mathbf{Z}}'$	$\tilde{\mathbf{D}}$	$\mathbf{T}$
Wald	$\mathbf{Z}'\mathbf{M}_{1_N}$	$\mathbf{M}_{1_N}\mathbf{D}$	$\mathbf{M}_{1_N}$
	$\uparrow \mathbf{P}_X = \mathbf{P}_{1_N} \uparrow$	$\uparrow \mathbf{P}_X = \mathbf{P}_{1_N} \uparrow$	$\uparrow \mathbf{P}_X = \mathbf{P}_{1_N} \uparrow$
TSLS	$\mathbf{Z}'\mathbf{M}_X$	$\mathbf{M}_X\mathbf{D}$	$\mathbf{M}_X$
	$\uparrow \hat{\mathbf{Z}} = \mathbf{P}_X\mathbf{Z} \uparrow$	$\uparrow \hat{\mathbf{D}} = \mathbf{P}_X\mathbf{D} \uparrow$	$\uparrow \hat{\mathbf{Y}} = \mathbf{P}_X\mathbf{Y} \uparrow$
PLR-IV	$\hat{\mathbf{R}}'$	$\hat{\mathbf{V}}$	$(\mathbf{I}_N - \mathbf{S})$
	$\uparrow \alpha^{if} = \mathbf{1}_N \uparrow$	$=$	$=$
IF	$\hat{\mathbf{R}}'diag(\alpha^{if}(\mathbf{x}))$	$\hat{\mathbf{V}}$	$(\mathbf{I}_N - \mathbf{S})$
	$\downarrow \hat{\mathbf{R}} = \hat{\mathbf{V}} \& \alpha^{if} = \alpha^{cf} \downarrow$	$=$	$=$
CF	$\hat{\mathbf{V}}'diag(\alpha^{cf}(\mathbf{x}))$	$\hat{\mathbf{V}}$	$(\mathbf{I}_N - \mathbf{S})$
	$\downarrow \alpha^{cf} = \mathbf{1}_N \downarrow$	$=$	$=$
PLR	$\hat{\mathbf{V}}'$	$\hat{\mathbf{V}}$	$(\mathbf{I}_N - \mathbf{S})$
	$\downarrow \hat{\mathbf{D}} = \mathbf{P}_X\mathbf{D} \downarrow$	$\downarrow \hat{\mathbf{D}} = \mathbf{P}_X\mathbf{D} \downarrow$	$\downarrow \hat{\mathbf{Y}} = \mathbf{P}_X\mathbf{Y} \downarrow$
OLS	$\mathbf{D}'\mathbf{M}_X$	$\mathbf{M}_X\mathbf{D}$	$\mathbf{M}_X$
	$\downarrow \mathbf{P}_X = \mathbf{P}_{1_N} \downarrow$	$\downarrow \mathbf{P}_X = \mathbf{P}_{1_N} \downarrow$	$\downarrow \mathbf{P}_X = \mathbf{P}_{1_N} \downarrow$
DiM	$\mathbf{D}'\mathbf{M}_{1_N}$	$\mathbf{M}_{1_N}\mathbf{D}$	$\mathbf{M}_{1_N}$

*Note:* Starting from the darkest row and following the arrows, the table shows how estimators follow as special cases by imposing restrictions in the white rows.

*Computational remark:* The original implementation in the `grf` package applies a constant in the weighted residual-on-residual regression. This complicates notation but Appendix A.3.1 provides the details how numerical equivalence between original output of `grf` and the weighted representation is obtained in the `OutcomeWeights` package.

### 3.2.2 Augmented inverse probability weighting and its special cases

Augmented inverse probability weighting (AIPW) is developed in a series of papers (e.g. [Robins, Rotnitzky, & Zhao, 1994, 1995](#); [Rotnitzky, Robins, & Scharfstein, 1998](#); [Chernozhukov et al., 2018](#)). AIPW is a PIVE with empirical moment condition

$$\mathbb{E}_N \left[ \left\{ \underbrace{\hat{Y}_{1,i}^d - \hat{Y}_{0,i}^d + \lambda_{1,i}^{ipw}(Y_i - \hat{Y}_{1,i}^d) - \lambda_{0,i}^{ipw}(Y_i - \hat{Y}_{0,i}^d)}_{=: \hat{Y}_i^{aipw}} - \hat{\tau}^{aipw} \underbrace{1}_{=: \hat{D}_i^{aipw}} \right\} \underbrace{1}_{=: \hat{Z}_i^{aipw}} \right] = 0 \quad (14)$$

and vector form

$$\hat{\tau}^{aipw} = (\mathbf{1}'_N \mathbf{1}_N)^{-1} \mathbf{1}'_N [\hat{\mathbf{Y}}_1^d - \hat{\mathbf{Y}}_0^d + \text{diag}(\boldsymbol{\lambda}_1^{ipw})(\mathbf{Y} - \hat{\mathbf{Y}}_1^d) - \text{diag}(\boldsymbol{\lambda}_0^{ipw})(\mathbf{Y} - \hat{\mathbf{Y}}_0^d)]. \quad (15)$$

The next step is to provide the transformation matrix. This is possible under Condition 1b that the outcome predictions are obtained by smoothers such that  $\mathbf{S}_d^d \mathbf{Y} = \hat{\mathbf{Y}}_d^d$ . Plugging this into (15) and rearranging delivers the transformation matrix

$$\hat{\tau}^{aipw} = N^{-1} \mathbf{1}'_N \underbrace{[\mathbf{S}_1^d - \mathbf{S}_0^d + \text{diag}(\boldsymbol{\lambda}_1^{ipw})(\mathbf{I}_N - \mathbf{S}_1^d) - \text{diag}(\boldsymbol{\lambda}_0^{ipw})(\mathbf{I}_N - \mathbf{S}_0^d)]}_{=: \mathbf{T}^{aipw}} \mathbf{Y} \quad (16)$$

and leads to the following result:<sup>5</sup>

**Corollary 2** (*outcome weights of AIPW*)

*Under Condition 1b such that the treatment specific outcome predictions can be written as  $\mathbf{S}_1^d \mathbf{Y} = \hat{\mathbf{Y}}_1^d$  and  $\mathbf{S}_0^d \mathbf{Y} = \hat{\mathbf{Y}}_0^d$ , the outcome weights of AIPW take the form*

$$\boldsymbol{\omega}^{aipw'} = N^{-1} \mathbf{1}'_N [\mathbf{S}_1^d - \mathbf{S}_0^d + \text{diag}(\boldsymbol{\lambda}_1^{ipw})(\mathbf{I}_N - \mathbf{S}_1^d) - \text{diag}(\boldsymbol{\lambda}_0^{ipw})(\mathbf{I}_N - \mathbf{S}_0^d)]. \quad (17)$$

Table 2 shows how RA can be obtained by setting all IPW weights to zero. IPW is recovered by setting all entries of the smoother matrices to zero.

<sup>5</sup>The AIPW implementation of the `grf` package uses an alternative moment condition. It is equivalent to (14) in expectation but uses different nuisance parameters and therefore differs numerically. However, also the outcome weights of this variant can be obtained as shown in [A.3.2](#).

Table 2: AIPW pseudo-variables and transformation matrices

	$\tilde{\mathbf{Z}}'$	$\tilde{\mathbf{D}}$	$\mathbf{T}$
RA	$\mathbf{1}'_N$	$\mathbf{1}_N$	$\mathbf{S}_1^d - \mathbf{S}_0^d$
	=	=	$\uparrow \boldsymbol{\lambda}_1^{ipw} = \boldsymbol{\lambda}_0^{ipw} = \mathbf{0}_N \uparrow$
AIPW	$\mathbf{1}'_N$	$\mathbf{1}_N$	$\mathbf{S}_1^d - \mathbf{S}_0^d + \text{diag}(\boldsymbol{\lambda}_1^{ipw})(\mathbf{I}_N - \mathbf{S}_1^d) - \text{diag}(\boldsymbol{\lambda}_0^{ipw})(\mathbf{I}_N - \mathbf{S}_0^d)$
	=	=	$\downarrow \mathbf{S}_1^d = \mathbf{S}_0^d = \mathbf{0}_{N \times N} \downarrow$
IPW	$\mathbf{1}'_N$	$\mathbf{1}_N$	$\text{diag}(\boldsymbol{\lambda}_1^{ipw} - \boldsymbol{\lambda}_0^{ipw})$

### 3.2.3 Wald-AIPW and its special cases

Tan (2006) propose an AIPW extension for the case of a binary instrument. This estimator has the same structure as the canonical Wald (1940) estimator but applies AIPW to estimate reduced form and first stage, respectively. Following Chernozhukov et al. (2018), the Wald-AIPW empirical moment condition in the form of Equation 2 reads

$$\mathbb{E}_N \left[ \left\{ \overbrace{\hat{Y}_{1,i}^z - \hat{Y}_{0,i}^z + \lambda_{1,i}^{ipw,z}(Y_i - \hat{Y}_{1,i}^z) - \lambda_{0,i}^{ipw,z}(Y - \hat{Y}_{0,i}^z)}^{\hat{Y}_i^{iv-aipw, z}} \right. \right. \quad (18)$$

$$\left. \left. - \hat{\tau}^{iv-aipw} \left( \underbrace{\hat{D}_{1,i}^z - \hat{D}_{0,i}^z + \lambda_{1,i}^{ipw,z}(D_i - \hat{D}_{1,i}^z) - \lambda_{1,i}^{ipw,z}(D_i - \hat{D}_{0,i}^z)}_{=: \tilde{D}_i^{iv-aipw}} \right) \right\} \underbrace{1}_{=: \tilde{Z}_i^{iv-aipw}} \right] = 0.$$

and in the form of Equation 4 becomes

$$\hat{\tau}^{iv-aipw} = (\mathbf{1}'_N [\hat{\mathbf{D}}_1^z - \hat{\mathbf{D}}_0^z + \text{diag}(\boldsymbol{\lambda}_1^{ipw,z})(\mathbf{D} - \hat{\mathbf{D}}_1^z) - \text{diag}(\boldsymbol{\lambda}_0^{ipw,z})(\mathbf{D} - \hat{\mathbf{D}}_0^z)])^{-1} \quad (19)$$

$$\times \mathbf{1}'_N [\hat{\mathbf{Y}}_1^z - \hat{\mathbf{Y}}_0^z + \text{diag}(\boldsymbol{\lambda}_1^{ipw,z})(\mathbf{Y} - \hat{\mathbf{Y}}_1^z) - \text{diag}(\boldsymbol{\lambda}_0^{ipw,z})(\mathbf{Y} - \hat{\mathbf{Y}}_0^z)].$$

Following similar steps as in Section 3.2.2 establishes another special case of Proposition 1:

**Corollary 3** (*outcome weights of Wald-AIPW*)

*Under Condition 1c such that the instrument specific outcome predictions can be written*

as  $\mathbf{S}_1^z \mathbf{Y} = \hat{\mathbf{Y}}_1^z$  and  $\mathbf{S}_0^z \mathbf{Y} = \hat{\mathbf{Y}}_0^z$ , the outcome weights of Wald-AIPW take the form

$$\begin{aligned} \omega^{iv-aipw'} &= (\mathbf{1}'_N \overbrace{[\hat{\mathbf{D}}_1^z - \hat{\mathbf{D}}_0^z + \text{diag}(\lambda_1^{ipw,z})(\mathbf{D} - \hat{\mathbf{D}}_1^z) - \text{diag}(\lambda_0^{ipw,z})(\mathbf{D} - \hat{\mathbf{D}}_0^z)]}^{\tilde{\mathbf{D}}^{iv-aipw} :=})^{-1} \\ &\quad \times \mathbf{1}'_N \underbrace{[\mathbf{S}_1^z - \mathbf{S}_0^z + \text{diag}(\lambda_1^{ipw,z})(\mathbf{I}_N - \mathbf{S}_1^z) - \text{diag}(\lambda_0^{ipw,z})(\mathbf{I}_N - \mathbf{S}_0^z)]}_{=: \mathbf{T}^{iv-aipw}}. \end{aligned} \quad (20)$$

Table 3 summarizes the involved manipulations to arrive at Wald-RA and -IPW applying similar transformations as for AIPW but for both reduced form and first stage.

Table 3: Wald-AIPW pseudo-variables and transformation matrices

	$\tilde{\mathbf{Z}}'$	$\tilde{\mathbf{D}}$	$\mathbf{T}$
Wald-RA	$\mathbf{1}'_N$	$\hat{\mathbf{D}}_1^z - \hat{\mathbf{D}}_0^z$	$\mathbf{S}_1^z - \mathbf{S}_0^z$
	=	$\uparrow \lambda_1^{ipw,z} = \lambda_0^{ipw,z} = \mathbf{0}_N \uparrow$	$\uparrow \lambda_1^{ipw,z} = \lambda_0^{ipw,z} = \mathbf{0}_N \uparrow$
Wald-AIPW	$\mathbf{1}'_N$	$\tilde{\mathbf{D}}^{iv-aipw}$ in (20)	$\tilde{\mathbf{T}}^{iv-aipw}$ in (20)
	=	$\downarrow \hat{\mathbf{D}}_1^z = \hat{\mathbf{D}}_0^z \downarrow = \mathbf{0}_N$	$\downarrow \mathbf{S}_1^z = \mathbf{S}_0^z = \mathbf{0}_{N \times N} \downarrow$
Wald-IPW	$\mathbf{1}'_N$	$\text{diag}(\lambda_1^{ipw,z} - \lambda_0^{ipw,z})$	$\text{diag}(\lambda_1^{ipw,z} - \lambda_0^{ipw,z})$

### 3.3 Consolidation

This section provides the first characterization of outcome weights for IF, CF, PLR(-IV), and (Wald-)AIPW (the supplementary [theory in action notebook](#) illustrates that the numerical equivalences hold also in practice). The results highlight that the availability of outcome weights depends on the estimator implementation. In particular, it requires to apply smoothers for the involved outcome regressions (C1). This excludes methods with non-differentiable objective functions and/or non-linear link functions for outcome prediction, such as Lasso, (penalized) logistic regression, or many neural network architectures. However, it is important to note that the choices for instrument and treatment nuisance parameters do not affect the availability of outcome weights.

Overall, the simple framework of Section 2 proves very handy for compactly deriving new functional forms of outcome weights and recovering known ones. This is interesting and practically useful in its own right, as the obtained weights can be applied in any established weight-based routine. Additionally, the framework provides a natural lens to investigate basic properties of the outcome weights as we pursue in the following.

## 4 Weights properties of pseudo-IV estimators

The results of the previous section enable users to *ex post* inspect whether outcome weights fulfill certain properties. For example, weights adding up to one for the treated (i.e.  $\sum_i \omega_i D_i = 1$ ) and to minus one for the untreated (i.e.  $\sum_i \omega_i (1 - D_i) = -1$ ) are often considered desirable because they guarantee certain in- and equivariances of estimators (Imbens & Rubin, 2015; Słoczyński et al., 2024). However, the PIVE framework also allows to analytically investigate the weights properties of estimator implementations. This is conceptually appealing and practically relevant because it permits *ex ante* control over weights properties. Specifically, we investigate under which conditions estimators fulfill one of the five weights properties collected in Table 4 spanned by the total, treated, and untreated weight sums, respectively (see Figure A.2 for a graphical illustration).

Table 4: Outcome weights classification

Weights property	$\sum_i \omega_i$	$\sum_i \omega_i D_i$	$\sum_i \omega_i (1 - D_i)$
fully-unnormalized	$\neq 0$	$\neq 1$	$\neq -1$
untreated-unnormalized	$\neq 0$	$= 1$	$\neq -1$
treated-unnormalized	$\neq 0$	$\neq 1$	$= -1$
scale-normalized	$= 0$	$= c \neq 1$	$= -c \neq -1$
fully-normalized	$= 0$	$= 1$	$= -1$

The literature documents examples for each class in Table 4.<sup>6</sup> Fully-unnormalized weights are associated with inverse probability weighting since Hájek (1971). (Un)treated-unnormalized weights recently appeared in estimators building on Abadie’s (2003)  $\kappa_0$  and  $\kappa_1$  where only one group shows weights adding up to (minus) one (Słoczyński et al., 2024). Scale-normalized weights are described by Słoczyński et al. (2023) in the context of covariate balancing propensity scores of Imai and Ratkovic (2014). Such estimators have treated (untreated) weights summing to (minus) the same non-one constant  $c$  and also appear prominently in the analysis of partially linear regression based estimators below. Fully-normalized weights are the norm (see overview in Imbens & Rubin, 2015, Ch. 19).

<sup>6</sup>The proposed class labels in Table 4 ensure that all three versions of unnormalized weights would also be labeled as unnormalized by Słoczyński et al. (2024). Although “normalized” is a loosely defined term, it seems reasonable in this context to use it for estimators whose outcome weights sum to zero, to remain consistent with previous work.



## 4.1 Weights properties in the PIVE framework - general

The outcome weights properties in Table 4 are determined by three weight sums. This motivates the following protocol to classify PIVE weights:

1. Calculate  $C := \sum_i \omega_i = \boldsymbol{\omega}' \mathbf{1}_N$ 
  - If  $C = 0 \Rightarrow$  normalized
  - If  $C \neq 0 \Rightarrow$  unnormalized
2. Calculate  $C_1 := \sum_i \omega_i D_i = \boldsymbol{\omega}' \mathbf{D}$ 
  - If  $C \neq 0$  and  $C_1 = 1 \Rightarrow$  untreated-unnormalized
  - If  $C = 0$  and  $C_1 \neq 1 \Rightarrow$  scale-normalized
  - If  $C = 0$  and  $C_1 = 1 \Rightarrow$  fully-normalized
3. If  $C \neq 0$  and  $C_1 \neq 1$ , calculate  $C_0 := \sum_i \omega_i (1 - D_i) = \boldsymbol{\omega}' (\mathbf{1}_N - \mathbf{D})$ 
  - If  $C_0 \neq -1 \Rightarrow$  fully-unnormalized
  - If  $C_0 = -1 \Rightarrow$  treated-unnormalized

Recall from Proposition 1 that PIVE weights take the form  $\boldsymbol{\omega}' = (\tilde{\mathbf{Z}}' \tilde{\mathbf{D}})^{-1} \tilde{\mathbf{Z}}' \mathbf{T}$ . Therefore classifying the weights properties of an estimator boils down to checking the first two or all of the following equations:

$$(\tilde{\mathbf{Z}}' \tilde{\mathbf{D}})^{-1} \tilde{\mathbf{Z}}' \mathbf{T} \mathbf{1}_N = 0 \tag{21}$$

$$(\tilde{\mathbf{Z}}' \tilde{\mathbf{D}})^{-1} \tilde{\mathbf{Z}}' \mathbf{T} \mathbf{D} = 1 \tag{22}$$

$$(\tilde{\mathbf{Z}}' \tilde{\mathbf{D}})^{-1} \tilde{\mathbf{Z}}' \mathbf{T} (\mathbf{1}_N - \mathbf{D}) = -1 \tag{23}$$

This implies that it suffices to investigate the following properties of the transformation matrix as shortcuts to classify the outcome weights:

1.  $\mathbf{T} \mathbf{1}_N = \mathbf{0}_N$  because it implies that Equation 21 holds
2.  $\mathbf{T} \mathbf{D} = \tilde{\mathbf{D}}$  because it implies that Equation 22 holds
3.  $\mathbf{T} (\mathbf{1}_N - \mathbf{D}) = -\tilde{\mathbf{D}}$  because it implies that Equation 23 holds

This shows how the PIVE structure offers substantial complexity reduction streamlining the derivations below to a large extent. It turns out that the weights properties are intimately tied to implementation choices as we first illustrate in the OLS example before moving to more involved cases.

*Example (OLS) continued:* Only one aspect of the implementation affects OLS weights properties in the sense of Table 4:

**Condition 2** (*covariate matrix with constant*)

The covariate matrix  $\mathbf{X}$  contains a column of ones, which is by convention the first column.

We can therefore write for a matrix with  $p$  covariates  $\mathbf{X}(1, \mathbf{0}'_p)' = \mathbf{1}_N$ .

Condition 2 is fulfilled in any reasonable application. However, making it explicit illustrates how implementation choices affect weights properties. We start by checking whether the weights sum to zero and use shortcut 1 focusing on the transformation matrix:

$$\mathbf{T}^{ols'} \mathbf{1}_N = \mathbf{M}_X \mathbf{1}_N = (\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')\mathbf{1}_N$$

$$\text{If C2} = \mathbf{I}_N \mathbf{1}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(1, \mathbf{0}'_p)' = \mathbf{1}_N - \mathbf{X}(1, \mathbf{0}'_p)' = \mathbf{0}_N \Rightarrow \text{normalized}$$

We conclude that OLS is always normalized if we include a constant. Next, we investigate whether weights of treated sum up to one via shortcut 2:

$$\mathbf{T}^{ols'} \mathbf{D} = \mathbf{M}_X \mathbf{D} = \tilde{\mathbf{D}}^{ols} \Rightarrow \text{untreated-unnormalized}$$

$$\Rightarrow \text{If C2} \Rightarrow \text{fully-normalized}$$

The transformation matrix applied to the treatment recovers the pseudo-treatment, which is sufficient for treated weights adding up to one. Curiously, this holds without further conditions implying that OLS is untreated-unnormalized even without a constant. Both results taken together recover a well-known fact that OLS is fully-normalized under Condition 2. Additionally, following the proposed framework step by step uncovers a nuisance regarding the case without a constant.

## 4.2 Weights properties in the PIVE framework - concrete

### 4.2.1 Implementation details

This section collectively introduces implementation details that become relevant in the later derivations. We start with a relatively mild condition:

**Condition 3** (*affine smoother matrix*)

*In addition to Condition 1, all rows of the smoother matrices add up to one:*

$$\mathbf{S}\mathbf{1}_N = \mathbf{S}_0^d\mathbf{1}_N = \mathbf{S}_1^d\mathbf{1}_N = \mathbf{S}_0^z\mathbf{1}_N = \mathbf{S}_1^z\mathbf{1}_N = \mathbf{1}_N$$

Most smoothers discussed in the literature fulfill this property. However, [Curth et al. \(2024\)](#) note that boosted trees can be an exception.

The next condition is relevant for the treatment group specific outcome nuisances:

**Condition 4** (*no smoothing between treatment groups*)

*In addition to Condition 1b, the treatment group specific predictions are formed using only observations of the respective group. This ensures that*

$$\mathbf{S}_1^d\mathbf{1}_N = \mathbf{S}_1^d\mathbf{D} \Rightarrow \mathbf{S}_1^d(\mathbf{1}_N - \mathbf{D}) = \mathbf{0}_N \quad (24)$$

$$\mathbf{S}_0^d\mathbf{1}_N = \mathbf{S}_0^d(\mathbf{1}_N - \mathbf{D}) \Rightarrow \mathbf{S}_0^d\mathbf{D} = \mathbf{0}_N \quad (25)$$

This condition is relevant for AIPW estimators and in line with standard implementations forming the group specific outcome models in the respective subgroups.

The next condition is less familiar but important for estimators based on partially linear regression and Wald-AIPW:

**Condition 5** (*outcome smoother matrix applied to treatment*)

*(C5a) The treatment predictions are formed using the outcome smoother matrix:*

$$\mathbf{S}\mathbf{D} = \hat{\mathbf{D}} \quad (26)$$

*(C5b) The treatment predictions in the different instrument groups are formed using the*

respective outcome smoother matrix:

$$\mathbf{S}_1^z \mathbf{D} = \hat{\mathbf{D}}_1^z \text{ and } \mathbf{S}_0^z \mathbf{D} = \hat{\mathbf{D}}_0^z \quad (27)$$

This goes against the idea of many flexible estimators to entertain different models for outcome and treatment predictions, respectively. Therefore, this condition is not in line with standard implementations.

The final condition is relevant for all estimators involving an inverse probability weighting component:

**Condition 6** (*normalized inverse probability weights*)

$$(C6a) \lambda_{0,i}^{norm} := \lambda_{0,i}^{ipw} / \mathbb{E}_N[\lambda_{0,i}^{ipw}] \text{ and } \lambda_{1,i}^{norm} := \lambda_{1,i}^{ipw} / \mathbb{E}_N[\lambda_{1,i}^{ipw}] \Rightarrow \mathbf{1}'_N \boldsymbol{\lambda}_1^{norm} = \mathbf{1}'_N \boldsymbol{\lambda}_0^{norm} = N$$

$$(C6b) \lambda_{0,i}^{norm,z} := \lambda_{0,i}^{ipw,z} / \mathbb{E}_N[\lambda_{0,i}^{ipw,z}] \text{ and } \lambda_{1,i}^{norm,z} := \lambda_{1,i}^{ipw,z} / \mathbb{E}_N[\lambda_{1,i}^{ipw,z}] \Rightarrow \mathbf{1}'_N \boldsymbol{\lambda}_1^{norm,z} = \mathbf{1}'_N \boldsymbol{\lambda}_0^{norm,z} = N$$

C6a is the standard [Hájek \(1971\)](#) normalization and usually recommended in applications ([Busso, DiNardo, & McCrary, 2014](#)). C6b is suggested by [Uysal \(2011\)](#) and recommended by [Słoczyński et al. \(2024\)](#).

#### 4.2.2 Weights properties of Instrumental Forest and its special cases

Without further conditions  $\boldsymbol{\omega}^{if}(\mathbf{x})$  in (11) is fully-unnormalized. In the following, we explore conditions leading to (fully-)normalized weights. First, we investigate how  $\mathbf{T}^{if} \mathbf{1}_N = \mathbf{0}_N$  could be obtained:

$$\mathbf{T}^{if} \mathbf{1}_N = (\mathbf{I}_N - \mathbf{S}) \mathbf{1}_N = \mathbf{1}_N - \mathbf{S} \mathbf{1}_N$$

$$\text{If } \mathbf{C3} = \mathbf{1}_N - \mathbf{1}_N = \mathbf{0}_N \Rightarrow \text{normalized}$$

This establishes that the standard implementation of IF in `grf` uses normalized weights because it applies the affine smoother random forest (C3) to estimate the outcome nuisance.

The next question is when treated weights sum to one. To this end, it is sufficient to understand when  $\mathbf{T}^{if} \mathbf{D} = \tilde{\mathbf{D}}^{if}$ :

$$\mathbf{T}^{if} \mathbf{D} = (\mathbf{I}_N - \mathbf{S}) \mathbf{D} = \mathbf{D} - \mathbf{S} \mathbf{D}$$

If C5a =  $\mathbf{D} - \hat{\mathbf{D}} = \hat{\mathbf{V}} = \tilde{\mathbf{D}}^{if} \Rightarrow$  untreated-unnormalized

$\Rightarrow$  If C3 & C5a  $\Rightarrow$  fully-normalized

The two results span different scenarios. The practically relevant one being that Condition 3 holds but Condition 5a does not because different treatment and outcome models are applied. This means that in practice IF weights are only scale-normalized but not fully-normalized. Only applying the same affine smoother matrix to predict outcome and treatment ensures fully-normalized weights.<sup>7</sup>

Recall from Table 1 that CF, PLR-IV and PLR use the same transformation matrix as IF. Consequently, they are also scale-normalized in standard applications. In contrast, OLS and TSLS apply the same projection matrix to form treatment and outcome predictions such that C5a holds by construction. Again the observations regarding OLS in the previous section immediately apply for TSLS because they share pseudo-treatment and transformation matrix. Also TSLS with a constant is fully-normalized and untreated-unnormalized without a constant.

For completeness observe that the difference in means estimator fulfils by construction both Conditions 3 and 5a, and is therefore always fully-normalized. An overview of conditions and weights properties is collected in Table 5 below.

### 4.2.3 Weights properties of AIPW and its special cases

First, we investigate under which conditions AIPW is normalized:

$$\mathbf{T}^{aipw} \mathbf{1}_N = \mathbf{S}_1^d \mathbf{1}_N - \mathbf{S}_0^d \mathbf{1}_N + \lambda_1^{ipw} - \text{diag}(\lambda_1^{ipw}) \mathbf{S}_1^d \mathbf{1}_N - \lambda_0^{ipw} + \text{diag}(\lambda_0^{ipw}) \mathbf{S}_0^d \mathbf{1}_N$$

$$\text{If C3} = \mathbf{1}_N - \mathbf{1}_N + \lambda_1^{ipw} - \lambda_1^{ipw} - \lambda_0^{ipw} + \lambda_0^{ipw} = \mathbf{0}_N \Rightarrow \text{normalized}$$

This result contains two surprising components. First, we did *not* apply normalized IPW weights (C6a) to achieve normalized AIPW. This means AIPW is self-normalizing once

---

<sup>7</sup>Curiously, applying the same non-affine smoother ensures that at least treated weights sum to one generalizing the observation regarding OLS without constant in the previous section.

affine smoothers are applied. Second, normalized IPW weights alone do not normalize AIPW weights as a similar simplification is not possible under C6a only.

The second step investigates when treated weights sum to one:

$$\mathbf{T}^{aipw} \mathbf{D} = \mathbf{S}_1^d \mathbf{D} - \mathbf{S}_0^d \mathbf{D} + \lambda_1^{ipw} - \text{diag}(\lambda_1^{ipw}) \mathbf{S}_1^d \mathbf{D} + \text{diag}(\lambda_0^{ipw}) \mathbf{S}_0^d \mathbf{D}$$

$$\text{If C4} = \mathbf{S}_1^d \mathbf{1}_N + \lambda_1^{ipw} - \text{diag}(\lambda_1^{ipw}) \mathbf{S}_1^d \mathbf{1}_N$$

$$\text{If C3 \& C4} = \mathbf{1}_N + \lambda_1^{ipw} - \lambda_1^{ipw} = \mathbf{1}_N = \tilde{\mathbf{D}}^{aipw} \Rightarrow \text{fully-normalized}$$

This means that standard implementations using affine smoothers to estimate outcome nuisances in the (un)treated groups separately are self-fully-normalizing regardless which IPW weights are applied. This implies that RA inherits weights properties from AIPW because it can be considered as applying IPW weights of zero (see Table 2). In contrast, IPW can be considered as applying smoother matrices of zeros. These uninformative smoother matrices by construction fulfill C4 but not C3 such that IPW weights are not (fully-)normalized. This recovers the well-known result of Hájek (1971) regarding IPW as a special case of AIPW. Obviously IPW with explicitly fully-normalized weights (under C6a) are fully-normalized.<sup>8</sup>

#### 4.2.4 Weights properties of Wald-AIPW and its special cases

We can not directly apply the results of Section 4.2.3 because the pseudo-outcome and -treatment differ. However, to show that the estimator is normalized if affine smoothers are applied for the outcome regressions requires only to change the superscripts:

$$\mathbf{T}^{iv-aipw} \mathbf{1}_N = \mathbf{S}_1^z \mathbf{1}_N - \mathbf{S}_0^z \mathbf{1}_N + \lambda_1^{ipw,z} - \text{diag}(\lambda_1^{ipw,z}) \mathbf{S}_1^z \mathbf{1}_N - \lambda_0^{ipw,z} + \text{diag}(\lambda_0^{ipw,z}) \mathbf{S}_0^z \mathbf{1}_N$$

$$\text{If C3} = \mathbf{1}_N - \mathbf{1}_N + \lambda_1^{ipw,z} - \lambda_1^{ipw,z} - \lambda_0^{ipw,z} + \lambda_0^{ipw,z} = \mathbf{0}_N \Rightarrow \text{normalized}$$

---

<sup>8</sup>To see this within the framework note that under C6a  $\mathbf{1}'_N \text{diag}(\lambda_1^{norm} - \lambda_0^{norm}) \mathbf{1}_N = N - N = 0$  and  $\text{diag}(\lambda_1^{norm} - \lambda_0^{norm}) \mathbf{D} = \mathbf{1}_N = \tilde{\mathbf{D}}^{aipw}$  establishing full-normalization.

However, the investigation of the sum of treated weights shows notable differences:

$$\begin{aligned}
T^{iv-aipw} D &= S_1^z D - S_0^z D + \text{diag}(\lambda_1^{ipw,z})(D - S_1^z D) - \text{diag}(\lambda_0^{ipw,z})(D - S_0^z D) \\
\text{If C5b} &= \hat{D}_1^z - \hat{D}_0^z + \text{diag}(\lambda_1^{ipw,z})(D - \hat{D}_1^z) - \text{diag}(\lambda_0^{ipw,z})(D - \hat{D}_0^z) \\
&= \tilde{D}^{iv-aipw} \Rightarrow \text{untreated-unnormalized} \\
&\Rightarrow \text{If C3 \& C5b} \Rightarrow \text{fully-normalized}
\end{aligned}$$

Wald-AIPW is therefore only scale-normalized unless we apply the outcome smoothers to also predict the treatments. This goes against the idea of using different models for each nuisance parameter. Unlike AIPW, Wald-AIPW is therefore not expected to be fully-normalized in standard applications. Another point worth noting is that separating the sample by instrument value when estimating outcome/treatment nuisances - an IV version of C4 - is not sufficient to achieve fully-normalized weights of Wald-AIPW.

Similar to the previous section Wald-RA inherits all properties from Wald-AIPW. However, Wald-IPW is always untreated-unnormalized because it can be considered as applying the same zero smoother matrix to outcome and treatment (C5b). Additionally normalizing the weights (C6b) makes Wald-IPW even fully-normalized.<sup>9</sup> This recovers observations regarding Wald-IPW by Słoczyński et al. (2024) within the framework of this paper. Also their findings regarding Abadie’s (2003)  $\kappa$  estimators can be obtained in the framework as shown in Appendix A.4.

### 4.3 Consolidation

Table 5 summarizes the sufficient conditions for closed-form and (fully-)normalized outcome weights.<sup>10</sup> Estimators with a check mark in the second column always have a weighted representation. Those are the estimators based on IPW and OLS where the weights are either obvious or at least well-studied (e.g. Imbens & Rubin, 1997, 2015; Imbens, 2015; Chattopadhyay & Zubizarreta, 2023). They are still included to demonstrate the

<sup>9</sup>This follows by considering the full numerator of the weight and not only the transformation matrix such that  $\mathbf{1}'_N \text{diag}(\lambda_1^{norm,z} - \lambda_0^{norm,z}) \mathbf{1}_N = N - N = 0$  due to Condition 6b.

<sup>10</sup>Table A.2 in the Appendix provides an extended table collecting which conditions are fulfilled by construction and including results for (un)treated-unnormalized weights for completeness. However, those are rather of academic value and we focus on the practically relevant cases in the main text.

Table 5: Conditions for closed-form and properties of outcome weights

Estimator	Closed-form	Normalized	Fully-normalized
Instrumental forest	C1a	C3	C3 & C5a
PLR-IV	C1a	C3	C3 & C5a
TSLS	✓	C2	C2
Wald	✓	✓	✓
Causal Forest	C1a	C3	C3 & C5a
PLR	C1a	C3	C3 & C5a
OLS	✓	C2	C2
DiM	✓	✓	✓
AIPW	C1b	C3	C3 & C4
RA	C1b	C3	C3 & C4
IPW	✓	C6a	C6a
Wald-AIPW	C1c	C3	C3 & C5b
Wald-RA	C1c	C3	C3 & C5b
Wald-IPW	✓	C6b	C6b

*Abbreviations:* DiM: difference in means; RA: regression adjustment; IPW: inverse probability weighting; AIPW: augmented IPW; OLS: ordinary least squares; PLR: partially linear regression; TSLS: two-stage least squares

generality of the framework but not to provide new insights. Those are obtained for more sophisticated outcome adaptive estimators for which weighted representations are not available in the literature.

The results collected in Table 5 highlight the crucial role of implementation details for availability and properties of outcome weights. First, column two documents that researchers can ensure that outcome weights can be accessed *ex post* by applying smoothers to form outcome predictions as shown in Section 3.2. Second, estimator specific implementation decisions *ex ante* determine certain weights properties. Columns three and four of Table 5 can serve as look-up table for researchers who want to ensure that a particular implementation of an estimator generates outcome weights of a desired class. They contain several surprising or at least undocumented results regarding six prominent DML and GRF instances:

1. PLR(-IV), causal/instrumental forests, and Wald-AIPW are not fully-normalized in standard implementations because they usually apply different treatment and



outcome models (C5 not fulfilled).

2. AIPW is fully-normalized in standard implementations because they usually apply affine smoothers and estimate treated and untreated outcomes separately (C3 & C4 fulfilled).

## 4.4 Empirical Monte Carlo illustration

This section runs an Empirical Monte Carlo Study (EMCS) to illustrate that most standard implementations of DML and GRF are not fully-normalized. EMCS take a real dataset and modify some components such that the ground truth is known in the semi-synthetic dataset (e.g. [Huber et al., 2013](#); [Wendling et al., 2018](#)). Here, we use the treatment, instrument, and covariates of the 401(k) data ([Chernozhukov, Hansen, & Spindler, 2016](#)) but with a noiseless outcome  $Y_i^* = 1 + D_i$ . This simulates the most powerful treatment leaving every untreated unit at one and shifting every treated unit to two. We expect estimators to estimate an effect of exactly one in this setting without outcome noise. However, only fully-normalized implementations are guaranteed to achieve this because for them,  $\omega'Y^* = \omega'(\mathbf{1}_N + \mathbf{D}) = 1$ .

This exercise is run with the `DoubleML` ([Bach, Kurz, Chernozhukov, Spindler, & Klaassen, 2024](#)) and the `grf` ([Tibshirani et al., 2024](#)) R packages applied to 100 bootstrap samples. The nuisance parameters in `DoubleML` are obtained using honest random forest (affine smoother) or `XGBoost` (non-affine smoother). Each function uses its default values. [Table 6](#) summarizes the ten implementations under consideration. The final column shows whether an implementation is fully-normalized according to the theoretical results in [Table 5](#) and therefore expected to find the “effect” of one.

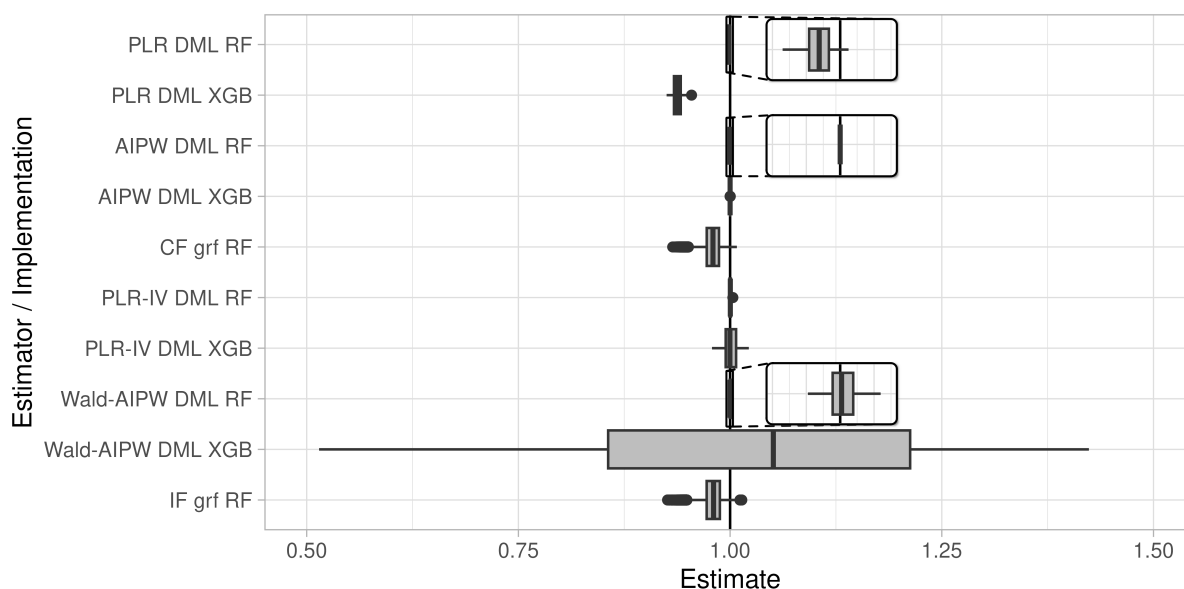
The theoretical predictions are confirmed in [Figure 1](#). The boxplots show that only AIPW with an affine smoother finds an effect of exactly one in all bootstrap samples. The other methods deviate from one to varying degrees. The `XGBoost` Wald-AIPW stands out in estimating effects between -16 and 55 (the graph is truncated). However, also causal/instrumental forests estimate heterogeneous effects between 0.93 and 1.01 although there is no heterogeneity to be found in the provided data. This illustrates the theoretical

Table 6: EMCS estimators and their labels

Label	Estimator	Package	Nuisance	Fully-normalized?
PLR DML RF	PLR	DoubleML	random forest	no b/c <a href="#">C5a</a>
PLR DML XGB	PLR	DoubleML	XGBoost	no b/c <a href="#">C3</a> & <a href="#">C5a</a>
AIPW DML RF	AIPW	DoubleML	random forest	yes
AIPW DML XGB	AIPW	DoubleML	XGBoost	no b/c <a href="#">C3</a>
CF grf RF	CF	grf	random forest	no b/c <a href="#">C5a</a>
PLR-IV DML RF	PLR-IV	DoubleML	random forest	no b/c <a href="#">C5a</a>
PLR-IV DML XGB	PLR-IV	DoubleML	XGBoost	no b/c <a href="#">C3</a> & <a href="#">C5a</a>
Wald-AIPW DML RF	Wald-AIPW	DoubleML	random forest	no b/c <a href="#">C5b</a>
Wald-AIPW DML XGB	Wald-AIPW	DoubleML	XGBoost	no b/c <a href="#">C3</a> & <a href="#">C5b</a>
IF grf RF	IF	grf	random forest	no b/c <a href="#">C5a</a>

*Notes:* The columns show (i) the labels used in Figure 1, (ii) which estimator defined in Section 3 is applied, (iii) the applied R package, (iv) the nuisance parameters, and (v) why the specific implementation are (not) expected to be fully-normalized.

Figure 1: Empirical Monte Carlo illustration



*Notes:* Boxplots show the results of 100 bootstraps of the 401(k) data ([Chernozhukov et al., 2016](#)) where the outcome is set to  $Y_i^* = 1 + D_i$ . The estimators are implemented using the default settings of the `DoubleML` and `grf` packages (see Table 6 for the labels). The causal/instrumental forest produces 9,915 estimates per replication such that their boxplots are based on  $\sim 1$  million estimates. The simulated effect is always one indicated by the solid line. The shadowed boxes in rows 1, 3 and 9 zoom into the range between 0.996 and 1.003. 49 outliers of Wald-AIPW DML XGB ranging from -16 to 55 are omitted. See [EMCS R notebook](#) for the code and more details.

findings even for DoubleML implementations where the extraction of the outcome weights is currently not possible because the required smoother matrices are not accessible.

## 5 Application: 401(k) covariate balancing

The novel outcome weights for DML and GRF can be used in established routines or to develop estimator-specific applications. We illustrate the former with covariate balancing, leaving the latter for future research. As in Section 4.4, we use the 401(k) data from [Chernozhukov et al. \(2016\)](#), but this time with the real outcome “net assets”.

### 5.1 Average effects

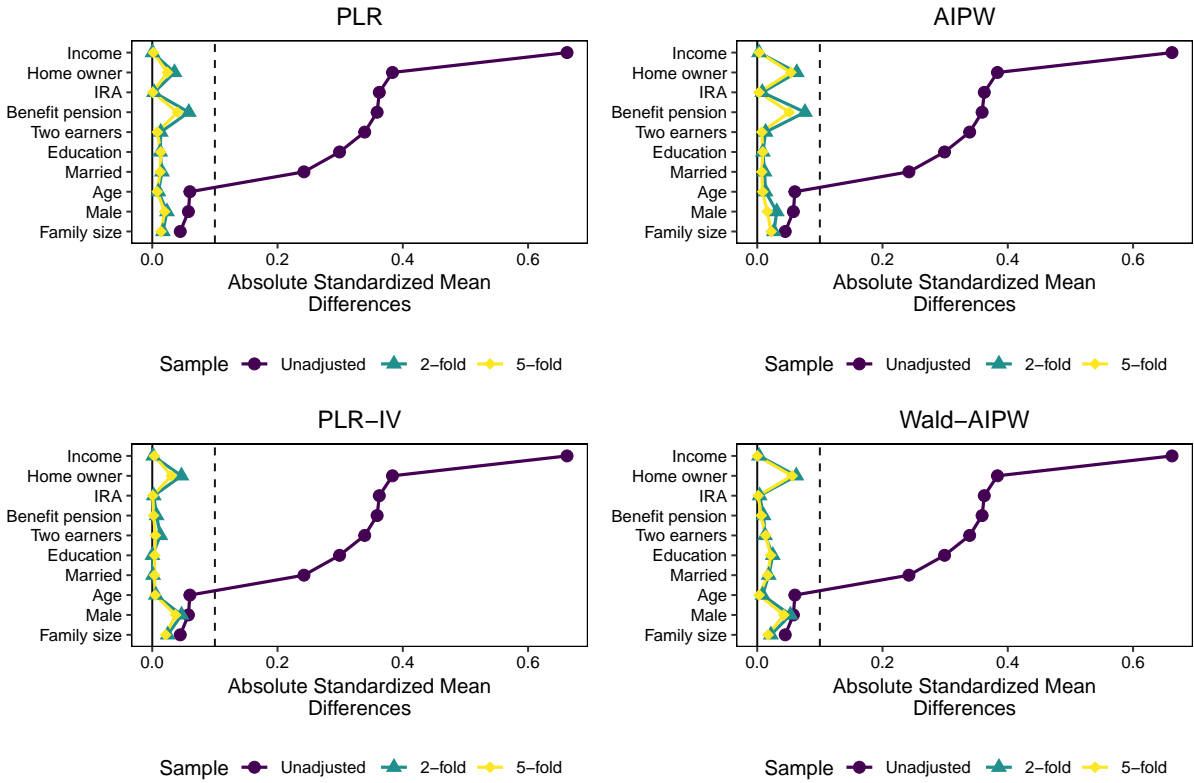
First, we investigate covariate balancing for DML estimated average effects. PLR(-IV) and (Wald-)AIPW are implemented using honest random forests with 2- and 5-fold cross-fitting. Figure 2 presents canonical balancing plots from the `cobalt` R package ([Greifer, 2024](#)) displaying absolute standardized mean differences (SMD). We observe that each method successfully balances the previously unbalanced covariates, in particular the income variable. Furthermore, cross-fitting with 5-folds achieves better covariate balancing compared to 2-folds. This demonstrates how DML outcome weights can be utilized in the design phase described by [Rubin \(2007\)](#), allowing researchers to commit to the preferred implementation before examining the results.

The supplementary [average effects R notebook](#) also provides point estimates and additional results, such as showing that 10 cross-fitting folds provide no further improvement over 5 folds and that the scale-normalized weights sum to values close to one (0.995 and closer).

### 5.2 Causal forest

Checks like those in Figure 2 are standard when estimating average effects. Similarly, we can assess covariate balancing for all 9,915 conditional average treatment effects (CATEs) produced by the `causal_forest()` function of the `grf` package. As an illustration, we investigate the importance of hyperparameter tuning for causal forests by comparing

Figure 2: Covariate balancing plots - average effects



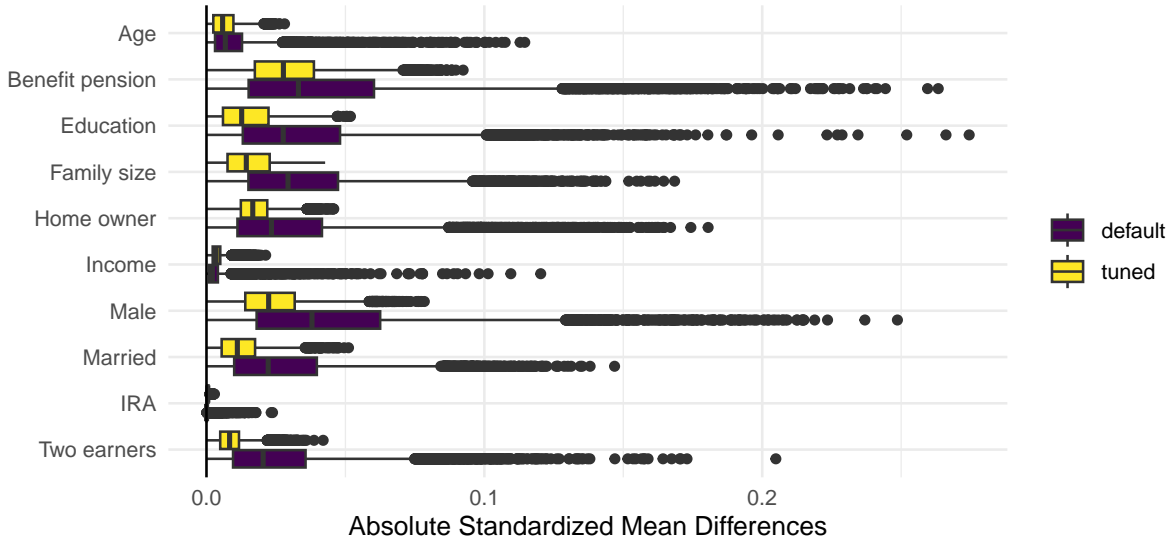
*Notes:* Each plot is created with the `love.plot()` function of the `cobalt` R package (Greifer, 2024) and the weights derived in Section 3.

the default implementation with `tune.parameters = "all"`. Figure 3 shows boxplots of absolute standardized mean differences (SMDs) for each CATE estimate. The results highlight that tuning the forests substantially improves covariate balancing in this application. The tuned version achieves absolute SMDs of 0.1 or lower, whereas the default settings frequently exceed this threshold, with some values even above 0.2. This highlights how standard diagnostics for average effects can also be applied to CATE estimates.

The supplementary [heterogeneous effects R notebook](#) reveals that the imbalances in the default forest coincide with implausible effect sizes ranging from  $-\$21k$  to  $\$78k$ , whereas the tuned forest yields more plausible estimates between  $\$8k$  and  $\$23k$ . This highlights the importance of parameter tuning for causal forests in this application. A similar pattern is observed for the instrumental forest, though with higher levels of  $|\text{SMD}|$ .

The supplementary notebook additionally examines descriptive statistics of the outcome weights multiplied by  $2D_i - 1$  to switch the sign of the untreated weights for better comparability. It documents that (i) both causal forests use negative weights, though to a

Figure 3: Covariate balancing plots - heterogeneous effects



*Notes:* Boxplots of absolute standardized mean differences for conditional average treatment effects estimated by causal forest using the default and tuned hyperparameters.

limited extent, (ii) instrumental forests assign substantial negative weights to never-takers, consistent with the outcome weights in [Imbens and Rubin \(1997\)](#) and the fact that the 401(k) setting has no always-takers by design, (iii) tuned forests use much smaller weights in absolute values, indicating more stable and reliable estimates, (iv) the sum of weights ranges from 0.98 to 1.02 for the default settings and from 0.995 to 1.005 for the tuned forest, making the tuned forest approximately fully-normalized in this application. Future work should explore whether this represents a general pattern.

## 6 Closing remarks

More estimators than previously noted can be expressed as weighted outcomes. The paper provides a general framework and derives novel weights for double machine learning and generalized random forest estimators. A key learning is that both availability and properties of the outcome weights depend on implementation choices and are therefore controlled by the user.

The paper focuses on providing general theoretical tools and standard illustrations. This acknowledges that access to their closed-form expressions is a prerequisite for developing new use cases or theoretical results for outcome weights. With the provided tools now

available, many follow-up questions arise for future research:

- Are there estimator specific use cases beyond the standard diagnostic tools?
- What are the closed-form expressions and properties of other PIVE outcome weights?
- Does the finding that several popular estimators do not use fully-normalized weights challenge the preference for such weights in the literature, or could explicitly normalizing the weights improve the finite sample performance of these estimators?
- Does the need to restrict outcome predictors to smoothers for access to outcome weights suggest a trade-off between interpretability and performance for outcome adaptive causal effect estimators?
- Do the provided outcome weights have implications for statistical inference, asymptotic properties, or double robustness robustness properties?

The investigation of the latter point most likely requires to restrict focus to analytically tractable smoothers in contrast to the generic smoothers allowed for in this paper. The fact that the smoothers and therefore the outcome weights may depend on the outcome pose non-trivial challenges. For example, it makes the outcome weights not compatible with approaches to use them for statistical inference as the existing approaches require outcome weights to be independent of the outcomes (e.g. [Imbens & Rubin, 2015](#), Ch. 19). Tailored sample splits as in [Lechner \(2018\)](#) could ensure the required independence but explorations along these lines are left for future work.

## References

- Abadie, A. (2003). Semiparametric instrumental variable estimation of treatment response models. *Journal of Econometrics*, *113*(2), 231–263.
- Abadie, A. (2021). Using synthetic controls: Feasibility, data requirements, and methodological aspects. *Journal of Economic Literature*, *59*(2), 391–425.
- Abadie, A., & Cattaneo, M. D. (2018). Econometric methods for program evaluation. *Annual Review of Economics*, *10*, 465–503.
- Angrist, J. D. (1998). Estimating the labor market impact of voluntary military service using social security data on military applicants. *Econometrica*, *66*(2), 249–288.
- Angrist, J. D., & Imbens, G. W. (1995). Two-stage least squares estimation of average causal effects in models with variable treatment intensity. *Journal of the American Statistical Association*, *90*(430), 431–442.
- Angrist, J. D., Imbens, G. W., & Rubin, D. B. (1996). Identification of Causal Effects Using Instrumental Variables. *Journal of the American Statistical Association*, *91*(434), 444–455.
- Angrist, J. D., & Krueger, A. B. (1999). Empirical strategies in labor economics. In O. C. Ashenfelter & D. Card (Eds.), *Handbook of labor economics* (Vol. 3A, pp. 1277–1366). Elsevier.
- Armstrong, T. B., & Kolesár, M. (2021). Finite-sample optimal estimation and inference on average treatment effects under unconfoundedness. *Econometrica*, *89*(3), 1141–1177.
- Aronow, P. M., & Samii, C. (2016). Does Regression Produce Representative Estimates of Causal Effects? *American Journal of Political Science*, *60*(1), 250–267.
- Athey, S., & Imbens, G. W. (2017). The state of applied econometrics: causality and policy evaluation. *Journal of Economic Perspectives*, *31*(2), 3–32.
- Athey, S., Tibshirani, J., & Wager, S. (2019). Generalized random forests. *Annals of Statistics*, *47*(2), 1148 - 1178.
- Athey, S., & Wager, S. (2019). Estimating treatment effects with causal forests: An application. *Observational Studies*, *5*(2), 37–51.
- Bach, P., Kurz, M. S., Chernozhukov, V., Spindler, M., & Klaassen, S. (2024). DoubleML: An object-oriented implementation of double machine learning in R. *Journal of Statistical Software*, *108*(3), 1–56.
- Ben-Michael, E., Feller, A., Hirshberg, D. A., & Zubizarreta, J. R. (2021). *The balancing act in causal inference*. Retrieved from <http://arxiv.org/abs/2110.14831>
- Blandhol, C., Bonney, J., Mogstad, M., & Torgovitsky, A. (2022). When is Tsls Actually Late? *SSRN Electronic Journal*. doi: 10.2139/ssrn.4021804
- Breitung, J., Bolwin, L., & Töns, J. (2024). *Alternative approaches for estimation and inference in synthetic control designs*.
- Bruns-Smith, D., Dukes, O., Feller, A., & Ogburn, E. L. (2023). *Augmented balancing weights as linear regression*. Retrieved from <http://arxiv.org/abs/2304.14545>

- Buja, A., Hastie, T., & Tibshirani, R. (1989). Linear smoothers and additive models. *The Annals of Statistics*, 17(2), 453–510.
- Busso, M., DiNardo, J., & McCrary, J. (2014). New evidence on the finite sample properties of propensity score reweighting and matching estimators. *Review of Economics and Statistics*, 96(5), 885–897.
- Chattopadhyay, A., & Zubizarreta, J. R. (2021). *On the implied weights of linear regression for causal inference*. Retrieved from <http://arxiv.org/abs/2104.06581v2>
- Chattopadhyay, A., & Zubizarreta, J. R. (2023). On the implied weights of linear regression for causal inference. *Biometrika*, 110(3), 615–629.
- Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., & Robins, J. (2018). Double/Debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1), C1-C68.
- Chernozhukov, V., Fernández-Val, I., Hahn, J., & Newey, W. (2013). Average and Quantile Effects in Nonseparable Panel Models. *Econometrica*, 81(2), 535–580.
- Chernozhukov, V., Hansen, C., Kallus, N., Spindler, M., & Syrgkanis, V. (2024). *Applied Causal Inference Powered by ML and AI*. forthcoming.
- Chernozhukov, V., Hansen, C., & Spindler, M. (2016). *High-Dimensional Metrics in R*. Retrieved from <http://arxiv.org/abs/1603.01700>
- Curth, A., Jeffares, A., & van der Schaar, M. (2023). A U-turn on double descent: Rethinking parameter counting in statistical learning. In *Thirty-seventh conference on neural information processing systems*.
- Curth, A., Jeffares, A., & van der Schaar, M. (2024). *Why do random forests work? Understanding tree ensembles as self-regularizing adaptive smoothers*. Retrieved from <https://arxiv.org/abs/2402.01502v1>
- de Chaisemartin, C., & D’Haultfœuille, X. (2023). Two-way fixed effects and differences-in-differences with heterogeneous treatment effects: a survey. *Econometrics Journal*, 26(3), C1-C30.
- Doudchenko, N., & Imbens, G. W. (2016). *Balancing, regression, difference-in-differences and synthetic control methods: A synthesis*. Retrieved from <https://arxiv.org/abs/1610.07748>
- Goldsmith-Pinkham, P., Hull, P., & Kolesár, M. (2021). *Contamination bias in linear regressions*. Retrieved from <http://arxiv.org/abs/2106.05024>
- Graham, B. S., Pinto, C. C. d. X., & Egel, D. (2012). Inverse probability tilting for moment condition models with missing data. *Review of Economic Studies*, 79(3), 1053–1079.
- Greifer, N. (2024). *cobalt: Covariate Balance Tables and Plots*. R package version 4.5.5.9000. Retrieved from <https://github.com/ngreifer/cobalt>
- Hainmueller, J. (2012). Entropy balancing for causal effects: A multivariate reweighting method to produce balanced samples in observational studies. *Political Analysis*, 20(1), 25–46.



- Hájek, J. (1971). Comment on “An essay on the logical foundations of survey sampling, part one”. In V. P. Godambe & D. A. Sprott (Eds.), *Foundations of statistical inference* (p. 236). Toronto: Holt, Rinehart and Winston.
- Hastie, T., & Tibshirani, R. (1990). Generalized additive models. *Monographs on Statistics and Applied Probability*, 43.
- Hazlett, C., & Shinkre, T. (2024). *Understanding and avoiding the "weights of regression": Heterogeneous effects, misspecification, and longstanding solutions*. Retrieved from <https://arxiv.org/abs/2403.03299>
- Heckman, J. J., & Vytlacil, E. (2005). Structural equations, treatment effects, and econometric policy evaluation. *Econometrica*, 73(3), 669–738.
- Heiler, P. (2022). Efficient covariate balancing for the local average treatment effect. *Journal of Business and Economic Statistics*, 40(4), 1569–1582.
- Horvitz, D. G., & Thompson, D. J. (1952). A generalization of sampling without replacement from a finite universe. *Journal of the American statistical Association*, 47(260), 663–685.
- Huber, M., Lechner, M., & Wunsch, C. (2013). The performance of estimators based on the propensity score. *Journal of Econometrics*, 175(1), 1–21.
- Humphreys, M. (2009). *Bounds on least squares estimates of causal effects in the presence of heterogeneous assignment probabilities*.
- Imai, K., & Ratkovic, M. (2014). Covariate balancing propensity score. *Journal of the Royal Statistical Society: Series B*, 76(1), 243–263.
- Imbens, G. W. (2004). Nonparametric estimation of average treatment effects under exogeneity: A review. *Review of Economics and Statistics*, 86(1), 4–29.
- Imbens, G. W. (2015). Matching methods in practice: Three examples. *Journal of Human Resources*, 50(2), 373–419. doi: 10.3368/jhr.50.2.373
- Imbens, G. W. (2024). Causal Inference in the Social Sciences. *Annual Review of Statistics and Its Application*, 11, 123–152.
- Imbens, G. W., & Angrist, J. D. (1994). Identification and estimation of local average treatment effects. *Econometrica*, 62(2), 467–475.
- Imbens, G. W., & Rubin, D. B. (1997). Estimating Outcome Distributions for Compliers in Instrumental Variables Models. *Review of Economic Studies*, 64(4), 555–574.
- Imbens, G. W., & Rubin, D. B. (2015). *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press.
- Imbens, G. W., & Wooldridge, J. M. (2009). Recent developments in the econometrics of program evaluation. *Journal of Economic Literature*, 47(1), 5–86.
- Jakiela, P. (2021). *Simple diagnostics for two-way fixed effects*. Retrieved from <https://arxiv.org/abs/2103.13229v1>
- Kallus, N. (2020). Generalized optimal matching methods for causal inference. *Journal of Machine Learning Research*, 21(62), 1–54.

- Khan, S., & Ugander, J. (2023). Adaptive normalization for IPW estimation. *Journal of Causal Inference*, 11(1).
- Kline, P. (2011). Oaxaca-Blinder as a reweighting estimator. *American Economic Review*, 101(3), 532–537.
- Knaus, M. C. (2021). A double machine learning approach to estimate the effects of musical practice on student’s skills. *Journal of the Royal Statistical Society. Series A: Statistics in Society*, 184(1), 282–300.
- Knaus, M. C. (2024). *OutcomeWeights*. R package version 0.1.0. Retrieved from <https://cran.r-project.org/web/packages/OutcomeWeights/index.html>
- Lechner, M. (2018). Modified causal forests for estimating heterogeneous causal effects. *arXiv:1812.09487*. Retrieved from <https://arxiv.org/abs/1812.09487>
- Lechner, M., & Strittmatter, A. (2019). Practical procedures to deal with common support problems in matching estimation. *Econometric Reviews*, 38(2), 193–207.
- Lee, D. S., & Lemieux, T. (2010). Regression discontinuity designs in economics. *Journal of Economic Literature*, 20(1), 281–355.
- Lin, Z., & Han, F. (2022). *On regression-adjusted imputation estimators of the average treatment effect*. Retrieved from <https://arxiv.org/abs/2212.05424v2>
- Robins, J. M., & Rotnitzky, A. (1995). Semiparametric efficiency in multivariate regression models with missing data. *Journal of the American Statistical Association*, 90(429), 122–129.
- Robins, J. M., Rotnitzky, A., & Zhao, L. P. (1994). Estimation of regression coefficients when some regressors are not always observed. *Journal of the American Statistical Association*, 89(427), 846–866.
- Robins, J. M., Rotnitzky, A., & Zhao, L. P. (1995). Analysis of semiparametric regression models for repeated outcomes in the presence of missing data. *Journal of the American Statistical Association*, 90(429), 106–121.
- Robinson, P. M. (1988). Root-n-consistent semiparametric regression. *Econometrica*, 56(4), 931.
- Rosenbaum, P. R., & Rubin, D. B. (1984). Reducing bias in observational studies using subclassification on the propensity score. *Journal of the American Statistical Association*, 79(387), 516–524.
- Rosenbaum, P. R., & Rubin, D. B. (1985). Constructing a control group using multivariate matched sampling methods that incorporate the propensity score. *The American Statistician*, 39(1), 33–38.
- Roth, J., Sant’Anna, P. H., Bilinski, A., & Poe, J. (2023). What’s trending in difference-in-differences? A synthesis of the recent econometrics literature. *Journal of Econometrics*, 235(2), 2218–2244.
- Rotnitzky, A., Robins, J. M., & Scharfstein, D. O. (1998). Semiparametric regression for repeated outcomes with nonignorable nonresponse. *Journal of the American Statistical Association*, 93(444), 1321–1339.

- Rubin, D. B. (2007). The design versus the analysis of observational studies for causal effects: Parallels with the design of randomized trials. *Statistics in Medicine*, 26(1), 20–36.
- Słoczyński, T. (2020). *When should we (not) interpret linear IV estimands as LATE?* Retrieved from <http://arxiv.org/abs/2011.06695>
- Słoczyński, T. (2022). Interpreting OLS estimands when treatment effects are heterogeneous: Smaller groups get larger weights. *Review of Economics and Statistics*, 104(3), 501–509.
- Słoczyński, T., Uysal, S. D., & Wooldridge, J. M. (2023). *Covariate balancing and the equivalence of weighting and doubly robust estimators of average treatment effects.* Retrieved from <https://arxiv.org/abs/2310.18563>
- Słoczyński, T., Uysal, S. D., & Wooldridge, J. M. (2024). Abadie’s Kappa and Weighting Estimators of the Local Average Treatment Effect. *Journal of Business and Economic Statistics*, forthcoming.
- Smith, J. A., & Todd, P. E. (2005). Does matching overcome LaLonde’s critique of nonexperimental estimators? *Journal of Econometrics*, 125(1-2), 305–353.
- Stone, C. J. (1977). Consistent nonparametric regression. *Annals of Statistics*, 5(4), 595–620.
- Tan, Z. (2006). Regression and weighting methods for causal inference using instrumental variables. *Journal of the American Statistical Association*, 101(476), 1607–1618.
- Tibshirani, J., Athey, S., Sverdrup, E., & Wager, S. (2024). *grf: Generalized Random Forests.*
- Uysal, D. (2011). Three Essays on Doubly Robust Estimation Methods. *PhD Dissertation, University of Konstanz.*
- Wald, A. (1940). The fitting of straight lines if both variables are subject to error. *The Annals of Mathematical Statistics*, 11(3), 284–300.
- Wendling, T., Jung, K., Callahan, A., Schuler, A., Shah, N. H., & Gallego, B. (2018). Comparing methods for estimation of heterogeneous treatment effects using observational data from health care databases. *Statistics in Medicine*, 37(23), 3309–3324.
- Zhao, Q. (2019). Covariate balancing propensity score by tailored loss functions. *Annals of Statistics*, 47(2), 965–993.
- Zubizarreta, J. R. (2015). Stable weights that balance covariates for estimation with incomplete outcome data. *Journal of the American Statistical Association*, 110(511), 910–922.

# A Supplementary Appendix

## A.1 Estimators under consideration

### A.1.1 Motivating target parameters

Let  $Y_i(1)$  and  $Y_i(0)$  be the potential outcomes under treatment and control, respectively. The paper is motivated by estimators of causal effects that aggregate the individual treatment effects  $Y_i(1) - Y_i(0)$  over different populations:

- $\mathbb{E}[Y_i(1) - Y_i(0)]$ , the average treatment effect (ATE)
- $\mathbb{E}[Y_i(1) - Y_i(0) \mid \mathbf{X}_i = \mathbf{x}]$ , the conditional ATE (CATE)
- $\mathbb{E}[Y_i(1) - Y_i(0) \mid \text{Complier}_i]$ , the local ATE (LATE), where  $\text{Complier}_i$  is the subgroup being shifted into treatment by a binary instrument ([Angrist, Imbens, & Rubin, 1996](#))
- $\mathbb{E}[Y_i(1) - Y_i(0) \mid \text{Complier}_i, \mathbf{X}_i = \mathbf{x}]$ , the conditional local ATE (CLATE)

Alternatively, we might impose a partially linear outcome model  $Y_i = \theta D_i + g(\mathbf{X}_i) + \epsilon_i$  and aim to estimate  $\theta$ .

Definition and identification of such parameters are discussed in detail in the literature and in textbooks. However, the numerical results provided in the main text also apply if the identifying assumptions do not hold, the target is explicitly non-causal, or the target is a different causal quantity.

### A.1.2 Estimators

Table [A.1](#) collects how the considered estimators differ in the aggregation level of the target effect (average or conditional effects), the research design in which they are usually applied (randomized controlled trials, unconfoundedness or instrumental variables), and regarding outcome modelling assumptions (none, partially linear, or linear models).

Table A.1: Overview of considered estimators

<b>Estimator</b>	<b>Aggregation level</b>	<b>Research design</b>	<b>Outcome model</b>	<b>Outcome weights in the literature</b>
DiM	Average	RCT	none	Imbens & Rubin (2015), Ch. 19.4
RA	Average	RCT/UC	none	-
IPW	Average	RCT/UC	none	Horvitz & Thompson (1952)
AIPW	Average	RCT/UC	none	Knaus (2021) (Post-Lasso) Chattopadhyay & Zubizarreta (2023) (OLS)
PLR	Average	RCT/UC	partially linear	-
OLS	Average	RCT/UC	linear	Chattopadhyay & Zubizarreta (2023)
Wald	Average	IV	none	Imbens & Rubin (1997)
Wald-RA	Average	IV	none	-
Wald-IPW	Average	IV	none	Abadie (2003)
Wald-AIPW	Average	IV	none	-
PLR-IV	Average	IV	partially linear	-
TSLs	Average	IV	linear	Chattopadhyay & Zubizarreta (2021)
Causal Forest	Conditional	RCT/UC	none	-
Instrumental Forest	Conditional	IV	none	-

*Abbreviations:* DiM: difference in means; RA: regression adjustment; IPW: inverse probability weighting; AIPW: augmented IPW; OLS: ordinary least squares; PLR: partially linear regression; TSLs: two-stage least squares; RCT: randomized controlled trials; UC: unconfoundedness; IV: instrumental variable

Figure A.1: Considered estimators and connections between them

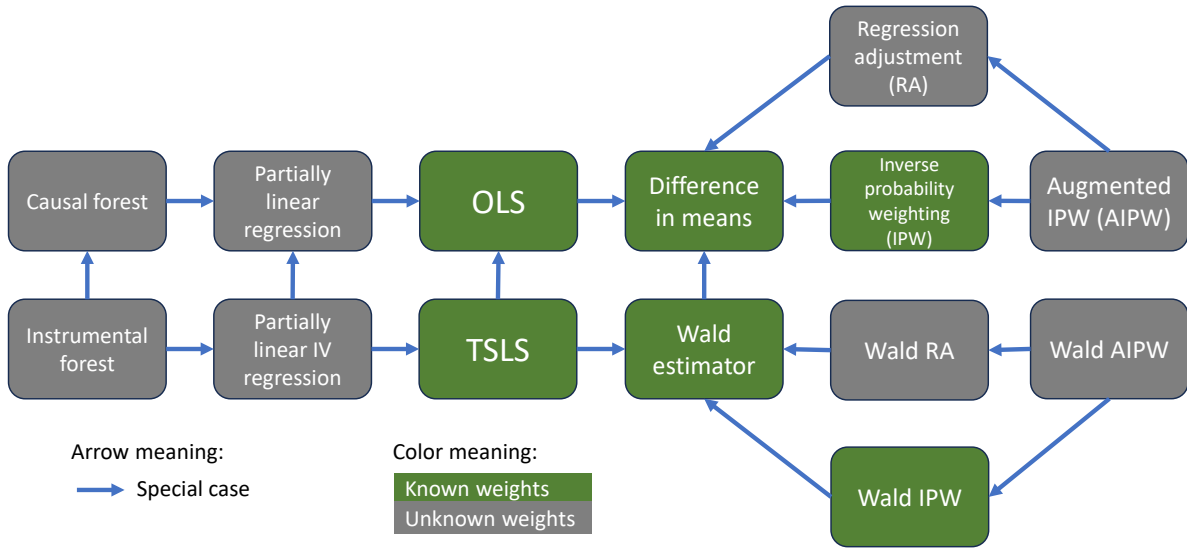


Figure A.1 additionally illustrates how estimators are connected and for which outcome weights are already known in the literature.

## A.2 More on smoothers

Many common regression estimators admit a representation as smoother. We distinguish three classes of smoothers:<sup>11</sup>

- $\mathbf{s}_i(X_i; \mathbf{X}, \mathbf{Y}, \epsilon_s)$  are *smoothers* that may depend on the outcome vector and on any type of randomness in building the prediction model, e.g. by inducing randomness in a random forest and/or by cross-validating the hyperparameters.
- $\mathbf{s}_i(X_i; \mathbf{X}, \mathbf{Y}, \epsilon_s) = \mathbf{s}_i(X_i; \mathbf{X}, \mathbf{Y}, \epsilon'_s) \forall \epsilon_s \neq \epsilon'_s$  are *deterministic smoothers* that do not depend on a random component, while still being outcome adaptive. One example would be (Post-)Lasso with data-driven penalty terms as implemented in the `hdm` R package of Chernozhukov et al. (2016).
- $\mathbf{s}_i(X_i; \mathbf{X}, \mathbf{Y}, \epsilon_s) = \mathbf{s}_i(X_i; \mathbf{X}, \mathbf{Y}', \epsilon'_s) \forall \epsilon_s \neq \epsilon'_s, \mathbf{Y} \neq \mathbf{Y}'$  are *linear smoothers* and neither depend on the outcome vector nor on a random component. OLS specified without using the data is a canonical linear smoother but also kernel and series

<sup>11</sup>See for a similar discussion the recent literature using smoothers to explain properties of machine learning methods (Curth, Jeffares, & van der Schaar, 2023; Curth et al., 2024).

regressions with fixed tuning parameter are linear smoothers (e.g. [Stone, 1977](#); [Buja, Hastie, & Tibshirani, 1989](#)).

The results in the main text merely require the existence of the smoother weights and do not depend on the smoother class. Therefore, we leave a more detailed discussion of the different classes for instances where the differences matter.

### A.3 grf package specific considerations

The main text ignores some complications arising in the R package `grf` implementing causal/instrumental forests and AIPW.

#### A.3.1 Causal forest

The first complication arises because the `grf` runs a weighted residual-on-residual regression *with constant*. The coefficient of this constant is typically not exactly zero because the weighted residuals are not guaranteed to sum to zero. Therefore, implementing Equation 13 will not exactly recover the package output. However, it can be achieved by defining the weighted least squares residual maker matrix  $\mathbf{M}_{\mathbf{1}_N}^\alpha := \mathbf{I}_N - \mathbf{1}_N(\mathbf{1}_N' \mathbf{diag}(\boldsymbol{\alpha}(\mathbf{x})) \mathbf{1}_N)^{-1} \mathbf{1}_N' \mathbf{diag}(\boldsymbol{\alpha}(\mathbf{x}))$  and using it in a modified version of (13) and for the causal forest as special case:

$$\boldsymbol{\omega}^{if}(\mathbf{x})' = (\hat{\mathbf{R}}' \mathbf{M}_{\mathbf{1}_N}^{\alpha^{if}} \mathbf{diag}(\boldsymbol{\alpha}^{if}(\mathbf{x})) \hat{\mathbf{V}})^{-1} \hat{\mathbf{R}}' \mathbf{M}_{\mathbf{1}_N}^{\alpha^{if}} \mathbf{diag}(\boldsymbol{\alpha}^{if}(\mathbf{x})) (\mathbf{I}_N - \mathbf{S}) \quad (\text{A.1})$$

$$\boldsymbol{\omega}^{cf}(\mathbf{x})' = (\hat{\mathbf{V}}' \mathbf{M}_{\mathbf{1}_N}^{\alpha^{cf}} \mathbf{diag}(\boldsymbol{\alpha}^{cf}(\mathbf{x})) \hat{\mathbf{V}})^{-1} \hat{\mathbf{V}}' \mathbf{M}_{\mathbf{1}_N}^{\alpha^{cf}} \mathbf{diag}(\boldsymbol{\alpha}^{cf}(\mathbf{x})) (\mathbf{I}_N - \mathbf{S}) \quad (\text{A.2})$$

This means we use a different pseudo-instrument compared to Equation 13 but leave pseudo-outcome and -treatment unchanged. Consequently, the conclusion in Section 4.2.2 that causal/instrumental forests are scale-normalized unless C5a is enforced remains valid because the pseudo-instrument does not affect the weights properties.

#### A.3.2 AIPW

The second complication arises when estimating the average treatment effect using AIPW with the `average_treatment_effect()` function. As described in [Athey and Wager \(2019\)](#)

Equation (8), the `grf` implementation applies an alternative representation of AIPW:

$$\mathbb{E}_N \left[ \underbrace{\hat{\tau}^{cf}(\mathbf{X}_i) + (\lambda_{1,i}^{ipw} - \lambda_{0,i}^{ipw})(Y_i - \hat{Y}_i - (D_i - \hat{D}_i)\hat{\tau}^{cf}(\mathbf{X}_i))}_{=:\hat{Y}_i^{aipw-grf}} - \hat{\tau}^{aipw-grf} \right] = 0 \quad (\text{A.3})$$

It uses the CATE estimates obtained by the causal forest  $\hat{\tau}^{cf}(\mathbf{X}_i)$  as nuisance parameter and not the two separate outcome regressions. To derive the outcome weights we store the CATEs of every observation in  $\hat{\boldsymbol{\tau}}^{cf} := (\hat{\tau}^{cf}(\mathbf{X}_1), \dots, \hat{\tau}^{cf}(\mathbf{X}_N))'$ . The solution of (A.3) in vector notation reads then

$$\hat{\boldsymbol{\tau}}^{aipw-grf} = N^{-1} \mathbf{1}'_N [\hat{\boldsymbol{\tau}}^{cf} + \text{diag}(\boldsymbol{\lambda}_1^{ipw} - \boldsymbol{\lambda}_0^{ipw})(\mathbf{Y} - \hat{\mathbf{Y}} - \text{diag}(\mathbf{D} - \hat{\mathbf{D}})\hat{\boldsymbol{\tau}}^{cf})]. \quad (\text{A.4})$$

We have established in Section 3.2.1 how to get the  $\mathbf{x}$ -specific causal forest weights and store them in a CATE smoother matrix  $\mathbf{S}^\tau := [\boldsymbol{\omega}^{cf}(\mathbf{X}_1) \dots \boldsymbol{\omega}^{cf}(\mathbf{X}_N)]'$  such that  $\hat{\boldsymbol{\tau}}^{cf} = \mathbf{S}^\tau \mathbf{Y}$ . The AIPW weights of the `grf` implementation are then

$$\boldsymbol{\omega}^{aipw-grf} = N^{-1} \mathbf{1}'_N [\mathbf{S}^\tau + \text{diag}(\boldsymbol{\lambda}_1^{ipw} - \boldsymbol{\lambda}_0^{ipw})(\mathbf{I}_N - \mathbf{S} - \text{diag}(\mathbf{D} - \hat{\mathbf{D}})\mathbf{S}^\tau)]. \quad (\text{A.5})$$

The next step is to investigate under which conditions `grf`-AIPW is normalized where we use  $\boldsymbol{\lambda}^{ipw} = \boldsymbol{\lambda}_1^{ipw} - \boldsymbol{\lambda}_0^{ipw}$  for compactness:

$$\begin{aligned} \mathbf{T}^{aipw-grf} \mathbf{1}_N &= [\mathbf{S}^\tau + \text{diag}(\boldsymbol{\lambda}^{ipw})(\mathbf{I}_N - \mathbf{S} - \text{diag}(\mathbf{D} - \hat{\mathbf{D}})\mathbf{S}^\tau)] \mathbf{1}_N \\ &= \mathbf{S}^\tau \mathbf{1}_N + \boldsymbol{\lambda}^{ipw} - \text{diag}(\boldsymbol{\lambda}^{ipw}) \mathbf{S} \mathbf{1}_N - \text{diag}(\boldsymbol{\lambda}^{ipw}) \text{diag}(\mathbf{D} - \hat{\mathbf{D}}) \mathbf{S}^\tau \mathbf{1}_N \end{aligned}$$

$$\text{If } \mathbf{C3} = \mathbf{0}_N + \boldsymbol{\lambda}^{ipw} - \boldsymbol{\lambda}^{ipw} - \text{diag}(\boldsymbol{\lambda}^{ipw}) \text{diag}(\mathbf{D} - \hat{\mathbf{D}}) \mathbf{0}_N = \mathbf{0}_N \Rightarrow \text{normalized}$$

because we have shown in Section 4.2.2 that causal forest is normalized under **C3** and therefore  $\mathbf{S}^\tau \mathbf{1}_N = \mathbf{0}_N$ . **C3** holds in the implementation of `grf` by default because it applies a random forest to estimate the outcome prediction. However, this is not enough



to guarantee fully-normalized weights:

$$\begin{aligned} \mathbf{T}^{aipw-grf} \mathbf{D} &= [\mathbf{S}^\tau + \text{diag}(\boldsymbol{\lambda}^{ipw})(\mathbf{I}_N - \mathbf{S} - \text{diag}(\mathbf{D} - \hat{\mathbf{D}})\mathbf{S}^\tau)] \mathbf{D} \\ &= \mathbf{S}^\tau \mathbf{D} + \boldsymbol{\lambda}_1^{ipw} - \text{diag}(\boldsymbol{\lambda}^{ipw}) \mathbf{S} \mathbf{D} - \text{diag}(\boldsymbol{\lambda}^{ipw}) \text{diag}(\mathbf{D} - \hat{\mathbf{D}}) \mathbf{S}^\tau \mathbf{D} \end{aligned}$$

$$\begin{aligned} \text{If C3 and C5a} &= \mathbf{1}_N + \boldsymbol{\lambda}_1^{ipw} - \text{diag}(\boldsymbol{\lambda}^{ipw}) \hat{\mathbf{D}} - \text{diag}(\boldsymbol{\lambda}^{ipw}) \text{diag}(\mathbf{D} - \hat{\mathbf{D}}) \mathbf{1}_N \\ &= \mathbf{1}_N + \boldsymbol{\lambda}_1^{ipw} - \text{diag}(\boldsymbol{\lambda}^{ipw}) \hat{\mathbf{D}} - \boldsymbol{\lambda}_1^{ipw} + \text{diag}(\boldsymbol{\lambda}^{ipw}) \hat{\mathbf{D}} = \mathbf{1}_N = \tilde{\mathbf{D}}^{aipw-grf} \\ &\Rightarrow \text{fully-normalized} \end{aligned}$$

This means that the AIPW estimator of `grf` is only scale-normalized because C5a is not fulfilled by default. This is in contrast to other implementations that are self-fully-normalized as discussed in Section 4.2.3.

#### A.4 Replicate Słoczyński et al. (2024) in the PIVE framework

Słoczyński et al. (2024) consider in total five estimators. Their estimators  $\hat{\tau}_t = \hat{\tau}_{a,1}$  correspond to the Wald-IPW without Condition 6b and  $\hat{\tau}_u$  to the Wald-IPW with normalized weights and are already discussed in Section 4.2.4. In addition the paper considers three estimators that require to define Abadie's (2003) kappas:

- $\kappa_0 := (1 - D_i) \frac{(1-Z_i)-(1-\hat{D}_i)}{\hat{D}_i(1-\hat{D}_i)}$
- $\kappa_1 := D_i \frac{Z_i-\hat{D}_i}{\hat{D}_i(1-\hat{D}_i)}$
- $\kappa := 1 - D_i \frac{1-Z_i}{1-\hat{D}_i} - (1 - D_i) \frac{Z_i(1-\hat{D}_i)}{\hat{D}_i}$

The three estimators are now presented in their vector form:

- $\hat{\tau}_a = (\mathbf{1}'_N \boldsymbol{\kappa})^{-1} \mathbf{1}'_N \text{diag}(\boldsymbol{\kappa}_1 - \boldsymbol{\kappa}_0) \mathbf{Y}$
- $\hat{\tau}_{a,0} = (\mathbf{1}'_N \boldsymbol{\kappa}_0)^{-1} \mathbf{1}'_N \text{diag}(\boldsymbol{\kappa}_1 - \boldsymbol{\kappa}_0) \mathbf{Y}$
- $\hat{\tau}_{a,10} = (\mathbf{1}'_N \mathbf{1}_N)^{-1} \mathbf{1}'_N \text{diag}(\boldsymbol{\kappa}_1 (\mathbf{1}'_N \boldsymbol{\kappa}_1)^{-1} N - \boldsymbol{\kappa}_0 (\mathbf{1}'_N \boldsymbol{\kappa}_0)^{-1} N) \mathbf{Y}$

Now we can apply the same strategies as in Section 4.2 to replicate that  $\hat{\tau}_a$  is fully-unnormalized,  $\hat{\tau}_{a,0}$  is treated-unnormalized, and  $\hat{\tau}_{a,10}$  is fully-normalized.

We note that  $\hat{\tau}_a$  and  $\hat{\tau}_{a,0}$  share the same transformation matrix and that the first two steps following shortcuts 1 and 2 are identical:

- $diag(\boldsymbol{\kappa}_1 - \boldsymbol{\kappa}_0)\mathbf{1}_N \neq \mathbf{0}_N \Rightarrow$  unnormalized
- $diag(\boldsymbol{\kappa}_1 - \boldsymbol{\kappa}_0)\mathbf{D} = \boldsymbol{\kappa}_1$ , which is not equal to the pseudo-treatments  $\tilde{\mathbf{D}}^{\hat{\tau}_a} = \boldsymbol{\kappa}$  or  $\tilde{\mathbf{D}}^{\hat{\tau}_{a,0}} = \boldsymbol{\kappa}_0$ , respectively  $\Rightarrow$  not untreated-unnormalized

To finish the characterization, we check whether the transformation matrix applied to the untreated produces minus the pseudo-treatment (shortcut 3):

$$\begin{aligned} diag(\boldsymbol{\kappa}_1 - \boldsymbol{\kappa}_0)(\mathbf{1}_N - \mathbf{D}) &= -\boldsymbol{\kappa}_0 \\ \Rightarrow \mathbf{T}^{\hat{\tau}_a}(\mathbf{1}_N - \mathbf{D}) &\neq -\tilde{\mathbf{D}}^{\hat{\tau}_a} \\ \Rightarrow \mathbf{T}^{\hat{\tau}_{a,0}}(\mathbf{1}_N - \mathbf{D}) &= -\tilde{\mathbf{D}}^{\hat{\tau}_{a,0}} \end{aligned}$$

We conclude that  $\hat{\tau}_a$  is fully unnormalized, but that  $\hat{\tau}_{a,0}$  is at least treated-unnormalized with untreated weights summing to minus one in line with [Słoczyński et al. \(2024\)](#).

Finally, we investigate  $\hat{\tau}_{a,10}$ . Here, it does not suffice to focus on the transformation matrix but we have to consider the full numerator to check whether weights sum to zero

$$\begin{aligned} \mathbf{1}'_N diag(\boldsymbol{\kappa}_1(\mathbf{1}'_N \boldsymbol{\kappa}_1)^{-1}N - \boldsymbol{\kappa}_0(\mathbf{1}'_N \boldsymbol{\kappa}_0)^{-1}N)\mathbf{1}_N &= \mathbf{1}'_N (\boldsymbol{\kappa}_1(\mathbf{1}'_N \boldsymbol{\kappa}_1)^{-1}N - \boldsymbol{\kappa}_0(\mathbf{1}'_N \boldsymbol{\kappa}_0)^{-1}N) \\ &= \mathbf{1}'_N \boldsymbol{\kappa}_1(\mathbf{1}'_N \boldsymbol{\kappa}_1)^{-1}N - \mathbf{1}'_N \boldsymbol{\kappa}_0(\mathbf{1}'_N \boldsymbol{\kappa}_0)^{-1}N \\ &= N - N = 0 \Rightarrow \text{normalized} \end{aligned}$$

and even the full expression to see that they are in addition fully-normalized:

$$\begin{aligned} (\mathbf{1}'_N \mathbf{1}_N)^{-1} \mathbf{1}'_N diag(\boldsymbol{\kappa}_1(\mathbf{1}'_N \boldsymbol{\kappa}_1)^{-1}N - \boldsymbol{\kappa}_0(\mathbf{1}'_N \boldsymbol{\kappa}_0)^{-1}N)\mathbf{D} &= N^{-1} \mathbf{1}'_N \boldsymbol{\kappa}_1(\mathbf{1}'_N \boldsymbol{\kappa}_1)^{-1}N \\ &= 1 \Rightarrow \text{fully-normalized} \end{aligned}$$

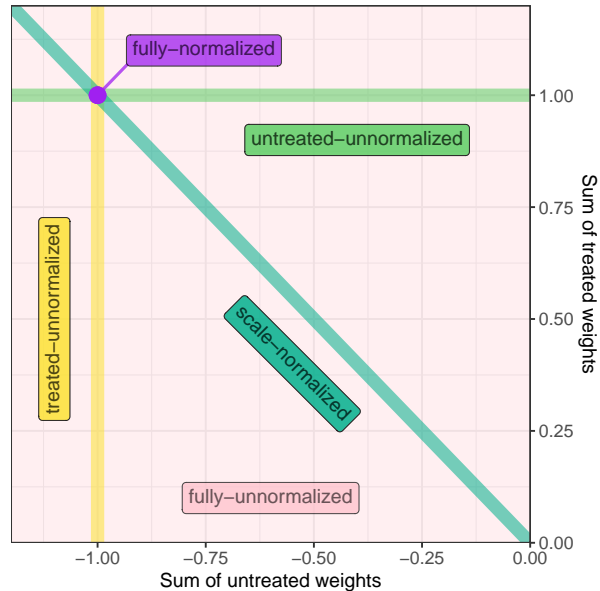
This final derivation highlights that the shortcuts 1-3 via the transformation matrix are only sufficient but not necessary to establish weights properties.

## A.5 More on outcome weights properties

### A.5.1 Graphical illustration of the classes

Figure A.2 illustrates the different classes formally defined in Table 4:

Figure A.2: Outcome weights classes



### A.5.2 Full results for weights properties

Table A.2 is an expanded version of Table 5 adding columns two, four and five. The second column stores which conditions are (not) fulfilled by construction for particular estimators, which influences the properties. The fourth and fifth column collect when estimators and (un)treated-normalized without being also normalized at the same time. These require rather artificial implementation decisions. For example in column four using affine smoothers only for the untreated but not for the treated outcome prediction ensures treated-unnormalized weights.

Table A.2: Conditions for closed-form and properties of outcome weights

Estimator	By construction	Closed-form	Treated-unnorm.	Untreated-unnorm.	Normalized	Fully-normalized
IF	-	C1a	-	C5a	C3	C3 & C5a
PLR-IV	-	C1a	-	C5a	C3	C3 & C5a
TSLS	C1a & C5a	✓	-	✓	C2	C2
Wald	C1a & C3 & C5b	✓	-	✓	✓	✓
CF	-	C1a	-	C5a	C3	C3 & C5a
PLR	-	C1a	-	C5a	C3	C3 & C5a
OLS	C1a & C5a	✓	-	✓	C2	C2
DiM	C1a & C3 & C5a	✓	-	-	✓	✓
AIPW	-	C1b	$S_0^d \mathbf{1}_N = \mathbf{1}_N$	$S_1^d \mathbf{1}_N = \mathbf{1}_N$	C3	C3 & C4
RA	-	C1b	$S_0^d \mathbf{1}_N = \mathbf{1}_N$	$S_1^d \mathbf{1}_N = \mathbf{1}_N$	C3	C3 & C4
IPW	C1b & C4 & <del>C3</del>	✓	with $\lambda_0^{norm}$	with $\lambda_1^{norm}$	C6a	C6a
Wald-AIPW	-	C1c	-	C5b	C3	C3 & C5b
Wald-RA	-	C1c	-	C5b	C3	C3 & C5b
Wald-IPW	C1c & C5b & <del>C3</del>	✓	-	✓	C6b	C6b

Notes: A “-” in columns 4 and 5 indicates that no condition of Section 4.2.1 leads to (un)treated-unnormalized weights without making them also normalized.