





## DoubleML: An Object-Oriented Implementation of Double Machine Learning in R

Philipp Bach   
University of Hamburg

Malte S. Kurz  
Technical University  
of Munich

Victor Chernozhukov   
Massachusetts Institute  
of Technology

Martin Spindler   
University of Hamburg

Sven Klaassen   
University of Hamburg

---

### Abstract

The R package **DoubleML** implements the double/debiased machine learning framework of Chernozhukov, Chetverikov, Demirer, Duflo, Hansen, Newey, and Robins (2018). It provides functionalities to estimate parameters in causal models based on machine learning methods. The double machine learning framework consists of three key ingredients: Neyman orthogonality, high-quality machine learning estimation and sample splitting. Estimation of nuisance components can be performed by various state-of-the-art machine learning methods that are available in the **mlr3** ecosystem. **DoubleML** makes it possible to perform inference in a variety of causal models, including partially linear and interactive regression models and their extensions to instrumental variable estimation. The object-oriented implementation of **DoubleML** enables a high flexibility for the model specification and makes it easily extendable. This paper serves as an introduction to the double machine learning framework and the R package **DoubleML**. In reproducible code examples with simulated and real data sets, we demonstrate how **DoubleML** users can perform valid inference based on machine learning methods.

*Keywords:* machine learning, causal inference, causal machine learning, R, **mlr3**, object orientation.

---

## 1. Introduction

Structural equation models provide a quintessential framework for conducting causal inference in statistics, econometrics, machine learning (ML), and other data sciences. The package **DoubleML** (Bach, Chernozhukov, Kurz, Spindler, and Sven 2024) for R (R Core Team 2023)

implements partially linear and interactive structural equation and treatment effect models with high-dimensional confounding variables as considered in [Chernozhukov \*et al.\* \(2018\)](#). Estimation and tuning of the machine learning models is based on the powerful functionalities provided by the **mlr3** package and the **mlr3** ecosystem ([Lang, Binder, Richter, Schratz, Pfisterer, Coors, Au, Casalicchio, Kotthoff, and Bischl 2019](#)). A key ingredient of double machine learning (DML) models are score functions identifying the estimates for the target parameter. These functions play an essential role for valid inference with machine learning methods because they have to satisfy a property called Neyman orthogonality. With the score functions as key elements, **DoubleML** implements double machine learning in a very general way using object orientation based on the **R6** package ([Chang 2021](#)). Currently, **DoubleML** implements the double / debiased machine learning framework as established in [Chernozhukov \*et al.\* \(2018\)](#) for

- partially linear regression models (PLR),
- partially linear instrumental variable regression models (PLIV),
- interactive regression models (IRM), and,
- interactive instrumental variable regression models (IIVM).

The object-oriented implementation of **DoubleML** is very flexible. The model classes ‘**DoubleMLPLR**’, ‘**DoubleMLPLIV**’, ‘**DoubleMLIRM**’ and ‘**DoubleIIVM**’ implement the estimation of the nuisance functions via machine learning methods and the computation of the Neyman-orthogonal score function. All other functionalities are implemented in the abstract base class ‘**DoubleML**’, including estimation of causal parameters, standard errors,  $t$  tests, confidence intervals, as well as valid simultaneous inference through adjustments of  $p$  values and estimation of joint confidence regions based on a multiplier bootstrap procedure. In combination with the estimation and tuning functionalities of **mlr3** and its ecosystem, this object-oriented implementation enables a high flexibility for the model specification in terms of

- the machine learning methods for estimation of the nuisance functions,
- the resampling schemes,
- the double machine learning algorithm, and,
- the Neyman-orthogonal score functions.

It further can be readily extended regarding

- new model classes that come with Neyman-orthogonal score functions being linear in the target parameter,
- alternative score functions via callables, and,
- customized resampling schemes.

Several other R packages for estimation of causal effects based on machine learning methods exist for R. The packages **grf** ([Tibshirani, Athey, and Wager 2023](#)) and **hdi** ([Dezeure, Bühlmann, Meier, and Meinshausen 2015](#)) implement alternative approaches to causal machine learning. **grf** implements generalized random forests ([Athey, Tibshirani, and Wager 2019](#)) and can be used for forest-based inference methods in different causal models including least-squares regression and estimation of treatment effects with and without instrumental variables. **hdi** can be used for inference in high-dimensional models with a focus on lasso-based estimation and methods for simultaneous inference.

An alternative approach, which was developed before the double machine learning framework, is the so-called targeted learning framework, and its software implementations and ecosystem (**tlverse**). For an overview and introduction to this approach and its implementations, we refer to the extensive **tlverse** handbook (Van der Laan, Coyle, Hejazi, Malenica, Phillips, and Hubbard 2022, <https://tlverse.org/tlverse-handbook>). Relevant R packages include **SuperLearner** (Polley, LeDell, Kennedy, and Van der Laan 2023) for flexible estimation using machine learning and **tmle** (Gruber and Van der Laan 2012) which implements estimation of causal parameters using targeted maximum likelihood estimation (TMLE). **sl3** (Coyle, Hejazi, Malenica, Phillips, and Sofrygin 2021) and **tmle3** (Coyle 2021) are recent extensions of the **tlverse** for object-oriented implementation of machine learning algorithms and a unified interface for TMLE.

Previous implementations that are more closely related to the double machine learning framework of Chernozhukov *et al.* (2018) have been provided by the R packages **hdm** (Chernozhukov, Hansen, and Spindler 2016), **dmlmt** (Knaus 2021), **causalDML** (Knaus 2022), **causalweight** (Bodory and Huber 2023) and **AIPW** (Zhong and Naimi 2021). **hdm** offers lasso-based inference methods in a variety of high-dimensional causal models, including estimation of (local) average treatment effects and linear (instrumental variable) regression. The underlying theoretical framework for valid post-selection and post-regularization inference can be considered as a special case of the more generic DML framework of Chernozhukov *et al.* (2018). Similarly, **dmlt** (Knaus 2021) provides methods of lasso-based inference on treatment effects of multi-valued treatment variables. Knaus (2022) focuses on a nonparametric treatment effect model, which is called the interactive regression model (IRM) in Chernozhukov *et al.* (2018) and also referred to as augmented inverse probability weighting. We adapted the term IRM from Chernozhukov *et al.* (2018) to denote this causal model and will refer to it accordingly in the following. The model is introduced in Section 3.3. Allowing for additional learners that include generalized random forests, ridge and random forests, **causalDML** (Knaus 2022) provides an implementation of the DML approach in an IRM in combination with recent methods for the analysis of heterogeneous treatment effects. Similarly, **causalweight** (Bodory and Huber 2023) focuses on the IRM and offers estimation methods for various causal quantities in this model as well as extensions thereof, including instrumental variable estimation, mediation analysis and sample selection approaches. In line with Chernozhukov *et al.* (2018), we denote the instrumental variable extension of the IRM in the following as IIVM. The model is introduced in Section 3.4. The R package **AIPW** (Zhong and Naimi 2021) implements estimation of average treatment effects of a binary treatment variable by augmented inverse probability weighting based on machine learning algorithms and integrates well with the **tlverse** ecosystem discussed later. In Python (Van Rossum *et al.* 2011), **EconML** (Battocchi, Dillon, Hei, Lewis, Oka, Oprescu, and Syrgkanis 2019) offers an implementation of several causal machine learning approaches. The package does not exclusively build on the double machine learning framework by Chernozhukov *et al.* (2018) and has a focus on heterogeneous effects.

In contrast to existing software packages, the R package **DoubleML** is intended to be a general implementation of the double machine learning approach of Chernozhukov *et al.* (2018). An introduction to the three key ingredients of the DML framework is provided in Section 4. The package can be used to perform inference in basically any causal model that can be characterized in terms of the formal framework of Chernozhukov *et al.* (2018). For example, it would be straightforward to extend **DoubleML** to mediation analysis (Farbmacher, Huber,

Laffers, Langen, and Spindler 2022), sample selection models (Bia, Huber, and Laffers 2020) or difference-in-differences (Chang 2020). As we will point out later, a key requirement for new model classes is a Neyman-orthogonal score. The object-oriented implementation makes the package easily extendable in terms of the supported causal models and other features of DML. By building on the **mlr3** ecosystem estimation can be based on a rich collection of powerful ML methods available in **mlr3** (Lang *et al.* 2019), **mlr3learners** (Lang, Au, Coors, and Schratz 2023a) and **mlr3extralearners** (Sonabend, Schratz, and Fischer 2023). The package **DoubleML** is available from the Comprehensive R Archive Network (CRAN) at <https://CRAN.R-project.org/package=DoubleML>.

We would like to mention that the R package **DoubleML** was developed together with a Python twin (Bach, Chernozhukov, Kurz, and Spindler 2022) that is based on **scikit-learn** (Pedregosa *et al.* 2011). The Python package is also available via GitHub (<https://github.com/DoubleML/doubleml-for-py>), the Python Package Index (PyPI, <https://pypi.org/project/DoubleML>), and conda-forge (<https://anaconda.org/conda-forge/doubleml>). Moreover, Kurz (2021) provides a serverless implementation of the Python module **DoubleML**.

The rest of the paper is structured as follows: In Section 2, we briefly demonstrate how to install the **DoubleML** package and give a short motivating example to illustrate the major idea behind the double machine learning approach. Section 3 introduces the main causal model classes implemented in **DoubleML**. Section 4 shortly summarizes the main ideas behind the double machine learning approach and reviews the key ingredients required for valid inference based on machine learning methods. Section 5 presents the main steps and algorithms of the double machine learning procedure for inference on one or multiple target parameters. Section 6 provides more detailed insights on the implemented classes and methods of **DoubleML**. Section 7 contains real-data and simulation examples for estimation of causal parameters using the **DoubleML** package. Additionally, this section provides a brief simulation study that illustrates the validity of the implemented methods in finite samples. Section 8 concludes the paper. The code output that has been suppressed in the main text and further information regarding the simulations are presented in the appendix. To make the code examples fully reproducible, the entire code is available in a supplementary zip file for this paper, as well as at <https://github.com/DoubleML/DoubleMLReplicationCode>. We would like to note that minor numerical differences might occur on other platforms when replicating code examples that involve random forest learners (see Appendix A for more information on the infrastructure used).

## 2. Getting started

### 2.1. Installation

The latest CRAN release of **DoubleML** can be installed using the command

```
R> install.packages("DoubleML")
```

Alternatively, the development version can be downloaded and installed from the GitHub (<https://github.com/DoubleML/doubleml-for-r>) repository using the command (previous installation of the **remotes** package is required, Csárdi, Hester, Wickham, Chang, Morgan, and Tenenbaum 2023)

```
R> remotes::install_github("DoubleML/doubleml-for-r")
```

Among others, **DoubleML** depends on the R package **R6** for object oriented implementation, **data.table** (Dowle and Srinivasan 2023) for the underlying data structure, as well as the packages **mlr3** (Lang *et al.* 2019), **mlr3learners** (Lang *et al.* 2023a) and **mlr3tuning** (Becker, Lang, Richter, Bischl, and Schalk 2023) for estimation of machine learning methods, model tuning and parameter handling. Moreover, the underlying packages of the machine learning methods that are called in **mlr3** or **mlr3learners** must be installed, for example the packages **glmnet** for lasso estimation (Friedman, Hastie, and Tibshirani 2010) or **ranger** (Wright and Ziegler 2017) for random forests.

Load the package after completed installation.

```
R> library("DoubleML")
```

## 2.2. A motivating example: Basics of double machine learning

In the following, we provide a brief summary of and motivation to double machine learning methods and show how the corresponding methods provided by the **DoubleML** package can be applied. The data generating process (DGP) is based on the introductory example in Chernozhukov *et al.* (2018). We consider a partially linear model: Our major interest is to estimate the causal parameter  $\theta$  in the following regression equation

$$y_i = \theta d_i + g_0(x_i) + \zeta_i, \quad \zeta_i \sim \mathcal{N}(0, 1),$$

with covariates  $x_i \sim \mathcal{N}(0, \Sigma)$ , where  $\Sigma$  is a matrix with entries  $\Sigma_{kj} = 0.7^{|j-k|}$ . In the following, the regression relationship between the treatment variable  $d_i$  and the covariates  $x_i$  will play an important role

$$d_i = m_0(x_i) + v_i, \quad v_i \sim \mathcal{N}(0, 1).$$

The nuisance functions  $m_0$  and  $g_0$  are given by

$$m_0(x_i) = x_{i,1} + \frac{1}{4} \frac{\exp(x_{i,3})}{1 + \exp(x_{i,3})},$$

$$g_0(x_i) = \frac{\exp(x_{i,1})}{1 + \exp(x_{i,1})} + \frac{1}{4} x_{i,3}.$$

We construct a setting with  $n = 500$  observations and  $p = 20$  explanatory variables to demonstrate the use of the estimators provided in **DoubleML**. Moreover, we set the true value of the parameter  $\theta$  to  $\theta = 0.5$ . The corresponding data generating process is implemented in the function `make_plr_CCDDHNR2018()`. We start by generating a realization of a data set as a ‘`data.table`’ object, which is subsequently used to create an instance of the data backend of class ‘`DoubleMLData`’.

```
R> library("DoubleML")
R> alpha <- 0.5
R> n_obs <- 500
R> n_vars <- 20
R> set.seed(1234)
R> data_plr <- make_plr_CCDDHNR2018(alpha = alpha, n_obs = n_obs,
+   dim_x = n_vars, return_type = "data.table")
```

The data backend implements the causal model: We specify that we perform inference on the effect of the treatment variable  $d_i$  on the dependent variable  $y_i$ .

```
R> obj_dml_data <- DoubleMLData$new(data_plr, y_col = "y", d_cols = "d")
```

In the next step, we choose the machine learning method as an object of class ‘Learner’ from **mlr3**, **mlr3learners** (Lang *et al.* 2023a) or **mlr3extralearners** (Sonabend *et al.* 2023). As we will point out later, we have to estimate two nuisance functions in order to perform valid inference in the partially linear regression model. Hence, we have to specify two learners. Moreover, we split the sample into two folds used for cross-fitting (`n_folds = 2`) in our illustrating examples for simplicity. Two-fold cross-fitting makes it necessary to estimate the ML models only twice, i.e., once per fold, which reduces the computational costs. In practice, it is generally recommended to choose a larger number of folds, cf. Remark 3. The default for the number of folds is `n_folds = 5`.

Load **mlr3** and **mlr3learners** packages and suppress output during estimation.

```
R> library("mlr3")
R> library("mlr3learners")
R> lgr::get_logger("mlr3")$set_threshold("warn")
```

Initialize a random forests learner with specified parameters.

```
R> ml_l <- lrn("regr.ranger", num.trees = 100, mtry = n_vars,
+   min.node.size = 2, max.depth = 5)
R> ml_m <- lrn("regr.ranger", num.trees = 100, mtry = n_vars,
+   min.node.size = 2, max.depth = 5)
R> ml_g <- lrn("regr.ranger", num.trees = 100, mtry = n_vars,
+   min.node.size = 2, max.depth = 5)
```

Initialize a causal model object, here a PLR.

```
R> doubleml_plr <- DoubleMLPLR$new(obj_dml_data,
+   ml_l, ml_m, ml_g, n_folds = 2, score = "IV-type")
```

To estimate the causal effect of variable  $d_i$  on  $y_i$ , we call the `fit()` method. To summarize the estimation output, we call the `summary()` method.

```
R> doubleml_plr$fit()
R> doubleml_plr$summary()
```

Estimates and significance testing of the effect of target variables

```
Estimate. Std. Error t value Pr(>|t|)
d 0.51457 0.04522 11.38 <2e-16 ***
---
```

```
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The output shows that the estimated coefficient is close to the true parameter  $\theta = 0.5$ . Moreover, there is evidence to reject the null hypotheses  $H_0 : \theta = 0$  at all common significance levels.

### 3. Key causal models

**DoubleML** provides estimation of causal effects in four different models: Partially linear regression models (PLR), partially linear instrumental variable regression models (PLIV), interactive regression models (IRM) and interactive instrumental variable regression models (IIVM). We will shortly introduce these models.

#### 3.1. Partially linear regression model

Partially linear regression models (PLR), which encompass the standard linear regression model, play an important role in data analysis (Robinson 1988). Partially linear regression models take the form

$$Y = D\theta_0 + g_0(X) + \zeta, \quad \mathbf{E}(\zeta \mid D, X) = 0, \quad (1)$$

$$D = m_0(X) + V, \quad \mathbf{E}(V \mid X) = 0, \quad (2)$$

where  $Y$  is the outcome variable and  $D$  is the policy variable of interest. The high-dimensional vector  $X = (X_1, \dots, X_p)$  consists of other confounding covariates, and  $\zeta$  and  $V$  are stochastic errors. Equation 1 is the equation of interest, and  $\theta_0$  is the main regression coefficient that we would like to infer. If  $D$  is conditionally exogenous (randomly assigned conditional on  $X$ ),  $\theta_0$  has the interpretation of a structural or causal parameter. The causal diagram supporting such interpretation is shown in Figure 1. The second equation keeps track of confounding, namely the dependence of  $D$  on covariates/controls. The characteristics  $X$  affect the policy variable  $D$  via the function  $m_0(X)$  and the outcome variable via the function  $g_0(X)$ . The partially linear model generalizes both linear regression models, where functions  $g_0$  and  $m_0$  are linear with respect to a collection of basis functions with respect to  $X$ , and approximately linear models.

An applied example from the economics literature is the analysis of the causal effect of 401(k) pension plans on employees' net financial assets by Poterba, Venti, and Wise (1994) and Poterba, Venti, and Wise (1995). In these studies, which are based on observational data, it is argued that eligibility for 401(k) pension plans can be assumed to be conditionally exogenous, once it is controlled for a set of confounders  $X$ , for example income. Following this argumentation and modelling approach, the estimate on  $\theta_0$  as obtained by a PLR can be interpreted as the average treatment effect of 401(k) eligibility on net financial assets. A reassessment and summary of the 401(k) example is available in Chernozhukov *et al.* (2018)

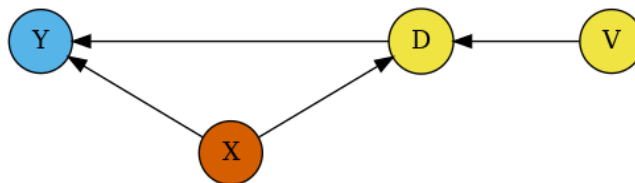


Figure 1: Causal diagram for PLR (Equation 1–2) and IRM (Equation 5–6) under conditional exogeneity. Note that the causal link between  $D$  and  $Y$  is one-directional. Identification of the causal effect is confounded by  $X$ , and identification is achieved via  $V$ , which captures variation in  $D$  that is independent of  $X$ . Methods to estimate the causal effect of  $D$  must therefore approximately remove the effect of high-dimensional  $X$  on  $Y$  and  $D$ .

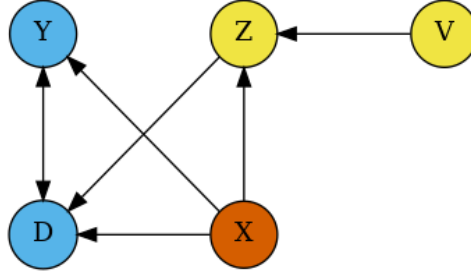


Figure 2: Causal diagram for PLIV (Equation 3–4) and IIVM (Equation 7–8) under conditional exogeneity of  $Z$ . Note that the causal link between  $D$  and  $Y$  is bi-directional, so an instrument  $Z$  is needed for identification. Identification is achieved via  $V$  that captures variation in  $Z$  that is independent of  $X$ . Equations 3 and 4 do not model the dependence between  $D$  and  $X$  and  $Z$ , though a necessary condition for identification is that  $Z$  and  $D$  are related after conditioning on  $X$ . Methods to estimate the causal effect of  $D$  must approximately remove the effect of high-dimensional  $X$  on  $Y$ ,  $D$ , and  $Z$ . Removing the confounding effect of  $X$  is done implicitly by the proposed procedure.

as well as on the **DoubleML** website ([https://docs.doubleml.org/stable/examples/R\\_double\\_ml\\_pension.html](https://docs.doubleml.org/stable/examples/R_double_ml_pension.html)).

### 3.2. Partially linear instrumental variable regression model

We next consider the partially linear instrumental variable regression model (PLIV)

$$Y - D\theta_0 = g_0(X) + \zeta, \quad \mathbb{E}(\zeta \mid Z, X) = 0, \quad (3)$$

$$Z = m_0(X) + V, \quad \mathbb{E}(V \mid X) = 0. \quad (4)$$

Note that this model is not a regression model unless  $Z = D$ . Model 3–4 is a canonical model in causal inference, going back to [Wright \(1928\)](#), with the modern difference being that  $g_0$  and  $m_0$  are nonlinear, potentially complicated functions of high-dimensional  $X$ . The idea of this model is that there is a structural or causal relation between  $Y$  and  $D$ , captured by  $\theta_0$ , and  $g_0(X) + \zeta$  is the stochastic error, partly explained by covariates  $X$ .  $V$  and  $\zeta$  are stochastic errors that are not explained by  $X$ . Since  $Y$  and  $D$  are jointly determined, we need an external factor, commonly referred to as an instrument,  $Z$ , to create exogenous variation in  $D$ . Note that  $Z$  should affect  $D$ . The  $X$  here serve again as confounding factors, so we can think of variation in  $Z$  as being exogenous only conditional on  $X$ .

A simple contextual example is from biostatistics ([Permutt and Hebel 1989](#)), where  $Y$  is a health outcome and  $D$  is an indicator of smoking. Thus,  $\theta_0$  captures the effect of smoking on health. Health outcome  $Y$  and smoking behavior  $D$  are treated as being jointly determined.  $X$  represents patient characteristics, and  $Z$  could be a doctor’s advice not to smoke (or another behavioral treatment) that may affect the outcome  $Y$  only through shifting the behavior  $D$ , conditional on characteristics  $X$ .

### 3.3. Interactive regression model

We consider estimation of average treatment effects when treatment effects are fully heterogeneous, i.e., the response curves under control and treatment can be different nonparametric



functions, and the treatment variable is binary,  $D \in \{0, 1\}$ . We consider vectors  $(Y, D, X)$  such that

$$Y = g_0(D, X) + U, \quad \mathbb{E}(U \mid X, D) = 0, \quad (5)$$

$$D = m_0(X) + V, \quad \mathbb{E}(V \mid X) = 0. \quad (6)$$

Since  $D$  is not additively separable, this model is more general than the partially linear model for the case of binary  $D$ . A common target parameter of interest in this model is the average treatment effect (ATE).

$$\theta_0 = \mathbb{E}[g_0(1, X) - g_0(0, X)].$$

Without unconfoundedness/conditional exogeneity, these quantities measure association, and could be referred to as average predictive effects (APE) and average predictive effect for the exposed (APEX). Inferential results for these objects would follow immediately from Theorem 1.

Another common target parameter is the average treatment effect for the treated (ATTE)

$$\theta_0 = \mathbb{E}[g_0(1, X) - g_0(0, X) \mid D = 1].$$

In business applications, the ATTE is often the main interest, as it captures the treatment effect for those who have been affected by the treatment. A difference of the ATTE from the ATE might arise if the characteristics of the treated individuals differ from those of the general population.

The confounding factors  $X$  affect the policy variable via the propensity score  $m_0(X)$  and the outcome variable via the function  $g_0(X)$ . Both of these functions are unknown and potentially complex, and we can employ ML methods to learn them.

Taking up the 401(k) example from Section 3.1, the general idea for identification of  $\theta_0$  using the IRM is similar. Once we are able to account for all confounding variables  $X$  in our analysis, we can consistently estimate the causal parameter  $\theta_0$ . A difference to the PLR refers to assumptions on the functional form of the main regression equation in 1 and 5, respectively. Whereas it is assumed that the effect of  $D$  on  $Y$  in the PLR model is additively separable, the IRM model comes with less restrictive assumptions. For example, it is possible that treatment effects are heterogeneous, i.e., vary across the population.

### 3.4. Interactive instrumental variable model

We consider estimation of the local average treatment effect (LATE) with a binary treatment variable  $D \in \{0, 1\}$ , and a binary instrument,  $Z \in \{0, 1\}$ . As before,  $Y$  denotes the outcome variable, and  $X$  is the vector of covariates. In a setting where unobserved factors drive the take-up of the treatment  $D$ , the average treatment effect is no longer identified. However, if a valid instrumental variable is available that changes individuals' decision to take up the treatment, it is possible to identify the LATE. The LATE measures the average causal effect for the subgroup of compliers, i.e., those individuals who receive the treatment only if the instrument takes value  $Z = 1$ . Hence, the LATE is of interest in many studies, where the treatment assignment cannot be assumed to be conditionally independent. For a more detailed treatment of the LATE and the key assumptions required for its identification, we would like to refer to [Imbens and Angrist \(1994\)](#), [Cunningham \(2021\)](#) and [Angrist and Pischke \(2009\)](#).

The structural equation for the IIVM is

$$Y = \ell_0(D, X) + \zeta, \quad \mathbb{E}(\zeta \mid Z, X) = 0, \quad (7)$$

$$Z = m_0(X) + V, \quad \mathbb{E}(V \mid X) = 0. \quad (8)$$

Consider the functions  $g_0$ ,  $r_0$ , and  $m_0$ , where  $g_0$  maps the support of  $(Z, X)$  to  $\mathbb{R}$  and  $r_0$  and  $m_0$  map the support of  $(Z, X)$  and  $X$  to  $(\epsilon, 1 - \epsilon)$  for some  $\epsilon \in (0, 1/2)$ , such that

$$Y = g_0(Z, X) + \nu, \quad \mathbb{E}(\nu \mid Z, X) = 0,$$

$$D = r_0(Z, X) + U, \quad \mathbb{E}(U \mid Z, X) = 0,$$

$$Z = m_0(X) + V, \quad \mathbb{E}(V \mid X) = 0.$$

We are interested in estimating

$$\theta_0 = \frac{\mathbb{E}[g_0(1, X)] - \mathbb{E}[g_0(0, X)]}{\mathbb{E}[r_0(1, X)] - \mathbb{E}[r_0(0, X)]}.$$

Under the well-known assumptions of [Imbens and Angrist \(1994\)](#),  $\theta_0$  is the LATE – the average causal effect for compliers, in other words, those observations that would have  $D = 1$  if  $Z$  were 1 and would have  $D = 0$  if  $Z$  were 0.

In the smoking example from [Section 3.2](#), the setting is similar to the section before, but now the binary treatment variable (“smoking”) is endogenous and is instrumented by a binary instrument variable  $Z$  (“doctor’s advice”). In this example, the group of compliers would comprise those individuals who quit smoking once their doctor advises them to do so and would otherwise continue to smoke. Similar to the comparison of the IRM model and the PLR model, the IIVM model does not impose the assumptions of linearity and additive separability that are maintained in the PLIV.

## 4. Basic idea and key ingredients of double machine learning

### 4.1. Basic idea behind double machine learning for the PLR model

Here we provide an intuitive discussion of how double machine learning works in the first model, the partially linear regression model. Naive application of machine learning methods directly to [Equations 1–2](#) may have a very high bias. Indeed, it can be shown that small biases in estimation of  $g_0$ , which are unavoidable in high-dimensional estimation, create a bias in the naive estimate of the main effect,  $\hat{\theta}_0^{\text{naive}}$ , which is sufficiently large to cause failure of conventional inference. The left panel in [Figure 3](#) illustrates this phenomenon. The histogram presents the empirical distribution of the studentized estimator,  $\hat{\theta}_0^{\text{naive}}$ , as obtained in 1000 independent repetitions of the data generating process presented in [Section 2.2](#). The functions  $g_0$  and  $m_0$  in the PLR model are estimated with random forest learners and corresponding predictions are then plugged into a non-orthogonal score function. The regularization performed by the random forest learner leads to a bias in estimation of  $g_0$  and  $m_0$ . Due to non-orthogonality of the score, this translates into a considerable bias of the main estimator  $\hat{\theta}_0^{\text{naive}}$ : The distribution of the studentized estimator  $\hat{\theta}_0^{\text{naive}}$  is shifted to the right of the origin and differs substantially from a normal distribution that would be obtained if the regularization bias was negligible as shown by the red curve.

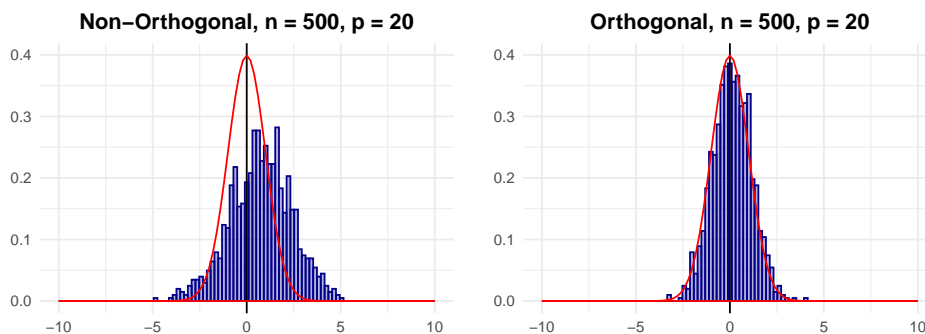


Figure 3: Performance of non-orthogonal and orthogonal estimators in a simulated data example. Left panel: Histogram of the studentized naive estimator  $\hat{\theta}_0^{\text{naive}}$ .  $\hat{\theta}_0^{\text{naive}}$  is based on estimation of  $g_0$  and  $m_0$  with random forests and a non-orthogonal score function. Data sets are simulated according to the data generating process in Section 2.2. Data generation and estimation are repeated 1000 times. Right panel: Histogram of the studentized DML estimator  $\tilde{\theta}_0$ .  $\tilde{\theta}_0$  is based on estimation of  $g_0$  and  $m_0$  with random forests and an orthogonal score function provided in Equation 11. Note that the simulated data sets and parameters of the random forest learners are identical to those underlying the left panel.

The PLR model above can be rewritten in the following residualized form

$$\begin{aligned} W &= V\theta_0 + \zeta, & \mathbb{E}(\zeta \mid D, X) &= 0, \\ W &= (Y - \ell_0(X)), & \ell_0(X) &= \mathbb{E}[Y \mid X], \\ V &= (D - m_0(X)), & m_0(X) &= \mathbb{E}[D \mid X]. \end{aligned}$$

The variables  $W$  and  $V$  represent original variables after taking out or *partially out* the effect of  $X$ . Note that  $\theta_0$  is identified from this equation if  $V$  has a non-zero variance.

Given identification, double machine learning for a PLR proceeds as follows

- (1) Estimate  $\ell_0$  and  $m_0$  by  $\hat{\ell}_0$  and  $\hat{m}_0$ , which amounts to solving the two problems of predicting  $Y$  and  $D$  using  $X$ , using any generic ML method, giving us estimated residuals

$$\hat{W} = Y - \hat{\ell}_0(X),$$

and

$$\hat{V} = D - \hat{m}_0(X).$$

The residuals should be of a cross-validated form, as explained below in Algorithm 1 or 2, to avoid biases from overfitting.

- (2) Estimate  $\theta_0$  by regressing the residual  $\hat{W}$  on  $\hat{V}$ . Use the conventional inference for this regression estimator, ignoring the estimation error in the residuals.

The reason we work with this residualized form is that it eliminates the bias arising from solving the prediction problems in stage (1). The estimates  $\hat{\ell}_0$  and  $\hat{m}_0$  carry a regularization bias due to having to solve prediction problems well in high-dimensions. However, the nature of the estimating equation for  $\theta_0$  are such that these biases are eliminated to the first order,

as explained below. This results in a high-quality low-bias estimator  $\tilde{\theta}_0$  of  $\theta_0$ , as illustrated in the right panel of Figure 3. The estimator is adaptive in the sense that the first stage estimation errors do not affect the second stage errors.

#### 4.2. Key ingredients of the double machine learning inference approach

Our goal is to construct high-quality point and interval estimators for  $\theta_0$  when  $X$  is high-dimensional and we employ machine learning methods to estimate the nuisance functions such as  $g_0$  and  $m_0$ . Example ML methods include lasso, random forests, boosted trees, deep neural networks, and ensembles or aggregated versions of these methods.

We shall use a method-of-moments estimator for  $\theta_0$  based upon the empirical analog of the moment condition

$$\mathbb{E}[\psi(W; \theta_0, \eta_0)] = 0, \quad (9)$$

where we call  $\psi$  the score function,  $W = (Y, D, X, Z)$ ,  $\theta_0$  is the parameter of interest, and  $\eta$  denotes nuisance functions with population value  $\eta_0$ .

*First key input: Neyman orthogonality*

The first key input of the inference procedure is using a score function  $\psi(W; \theta; \eta)$  that satisfies 9, with  $\theta_0$  being the unique solution, and that obeys the Neyman orthogonality condition

$$\partial_\eta \mathbb{E}[\psi(W; \theta_0, \eta)]|_{\eta=\eta_0} = 0. \quad (10)$$

Neyman orthogonality 10 ensures that the moment condition 9 used to identify and estimate  $\theta_0$  is insensitive to small perturbations of the nuisance function  $\eta$  around  $\eta_0$ . The derivative  $\partial_\eta$  denotes the pathwise (Gateaux) derivative operator.

In general, it is important to distinguish whether machine learning methods are used for prediction or in the context of statistical inference. An accurate prediction rule for the nuisance parameters  $\eta_0$  does not necessarily lead to a consistent estimator for the causal parameter  $\theta_0$ . Replacing the true value of  $\eta_0$  by an ML estimator  $\hat{\eta}_0$  likely introduces a bias, for example, due to heavy regularization in high-dimensional settings. If this bias is not taken into account, the estimator  $\hat{\theta}_0$  will generally be inconsistent and not have an asymptotically normal distribution. Using a Neyman-orthogonal score makes estimation of the causal parameter  $\theta_0$  robust against first order biases that arise from regularization. The Neyman orthogonality property is responsible for the adaptivity of the DML estimator – namely, the approximate distribution of  $\hat{\theta}_0$  will not depend on the fact that the estimate  $\hat{\eta}_0$  contains error, if the latter is mild. Other approaches, as targeted maximum likelihood and semiparametric sieves estimation recognize this as well. For a more detailed treatment of Neyman orthogonality we refer to Chernozhukov *et al.* (2018).

The right panel of Figure 3 presents the empirical distribution of the studentized DML estimator  $\tilde{\theta}_0$  that is based on an orthogonal score. Note that estimation is performed on the identical simulated data sets and with the same machine learning method as for the naive learner, which is displayed in the left panel. The histogram of the studentized estimator  $\tilde{\theta}_0$  illustrates the favorable performance of the double machine learning estimator, which is based on an orthogonal score: The DML estimator is robust to the bias that is generated by

regularization. The estimator is approximately unbiased, is concentrated around 0 and the distribution is well-approximated by the normal distribution.

PLR score: In the PLR model, we can employ two alternative score functions. We will shortly indicate the option for initialization of a model object in **DoubleML** to clarify how each score can be implemented. Using the option `score = "partialling out"` leads to estimation of the score function

$$\begin{aligned}\psi(W; \theta, \eta) &:= (Y - \ell(X) - \theta(D - m(X))) (D - m(X)), \\ \eta &= (\ell, m), \quad \eta_0 = (\ell_0, m_0),\end{aligned}\tag{11}$$

where  $W = (Y, D, X)$  and  $\ell$  and  $m$  are  $P$ -square-integrable functions mapping the support of  $X$  to  $\mathbb{R}$ , whose true values are given by

$$\ell_0(X) = \mathbb{E}[Y \mid X], \quad m_0(X) = \mathbb{E}[D \mid X].$$

Alternatively, it is possible to use the following score function for the PLR via the option `score = "IV-type"`

$$\psi(W; \theta, \eta) := (Y - D\theta - g(X)) (D - m(X)), \quad \eta = (g, m), \quad \eta_0 = (g_0, m_0),$$

with  $g$  and  $m$  being  $P$ -square-integrable functions mapping the support of  $X$  to  $\mathbb{R}$  with values given by

$$g_0 = \mathbb{E}[Y - D\theta_0 \mid X], \quad m_0(X) = \mathbb{E}[D \mid X].$$

The scores above are Neyman-orthogonal by elementary calculations. Now, it is possible to see the connections to the residualized system of equations presented in Section 4.1.

PLIV score: In the PLIV model, we can employ two alternative score functions. Using the option `score = "partialling out"` leads to estimation of the score function

$$\begin{aligned}\psi(W; \theta, \eta) &:= (Y - \ell(x) - \theta(D - r(X))) (Z - m(X)), \\ \eta &= (\ell, m, r), \quad \eta_0 = (\ell_0, m_0, r_0),\end{aligned}$$

where  $W = (Y, D, X, Z)$  and  $\ell$ ,  $m$ , and  $r$  are  $P$ -square integrable functions mapping the support of  $X$  to  $\mathbb{R}$ , whose true values are given by

$$\ell_0(X) = \mathbb{E}[Y \mid X], \quad r_0(X) = \mathbb{E}[D \mid X], \quad m_0(X) = \mathbb{E}[Z \mid X].$$

Alternatively, it is possible to use the following score function for the PLIV via the option `score = "IV-type"`

$$\psi(W; \theta, \eta) := (Y - D\theta - g(X)) (Z - m(X)), \quad \eta = (g, m), \quad \eta_0 = (g_0, m_0),$$

with  $g$  and  $m$  being  $P$ -square-integrable functions mapping the support of  $X$  to  $\mathbb{R}$  with values given by

$$g_0 = \mathbb{E}[Y - D\theta_0 \mid X], \quad m_0(X) = \mathbb{E}[Z \mid X].$$

IRM score: For estimation of the ATE parameter of the IRM model, we employ the score (`score = "ATE"`)

$$\begin{aligned}\psi(W; \theta, \eta) &:= (g(1, X) - g(0, X)) + \frac{D(Y - g(1, X))}{m(X)} - \frac{(1 - D)(Y - g(0, X))}{1 - m(X)} - \theta, \\ \eta &= (g, m), \quad \eta_0 = (g_0, m_0),\end{aligned}$$

where  $W = (Y, D, X)$  and  $g$  and  $m$  map the support of  $(D, X)$  to  $\mathbb{R}$  and the support of  $X$  to  $(\epsilon, 1 - \epsilon)$ , respectively, for some  $\epsilon \in (0, 1/2)$ , whose true values are given by

$$g_0(D, X) = \mathbb{E}[Y \mid D, X], \quad m_0(x) = \mathbb{P}[D = 1 \mid X].$$

This orthogonal score is based on the influence function for the mean for missing data from [Robins and Rotnitzky \(1995\)](#). For estimation of the ATTE parameter in the IRM, we use the score (`score = "ATTE"`)

$$\begin{aligned}\psi(W; \theta, \eta) &:= \frac{D(Y - g(0, X))}{p} - \frac{m(X)(1 - D)(Y - g(0, X))}{p(1 - m(x))} - \frac{D}{p}\theta, \\ \eta &= (g, m, p), \quad \eta_0 = (g_0, m_0, p_0),\end{aligned}$$

where  $p_0 = \mathbb{P}(D = 1)$ . Note that this score does not require estimating  $g_0(1, X)$ .

IIVM score: To estimate the LATE parameter in the IIVM, we will use the score (`score = "LATE"`)

$$\begin{aligned}\psi &:= g(1, X) - g(0, X) + \frac{Z(Y - g(1, X))}{m(X)} - \frac{(1 - Z)(Y - g(0, X))}{1 - m(X)} \\ &\quad - \left( r(1, x) - r(0, X) + \frac{Z(D - r(1, x))}{m(X)} - \frac{(1 - Z)(D - r(0, X))}{1 - m(X)} \right) \times \theta, \\ \eta &= (g, m, r), \quad \eta_0 = (g_0, m_0, r_0),\end{aligned}$$

where  $W = (Y, D, X, Z)$  and the nuisance parameter  $\eta = (g, m, r)$  consists of  $P$ -square integrable functions  $g$ ,  $m$ , and  $r$ , with  $g$  mapping the support of  $(Z, X)$  to  $\mathbb{R}$  and  $m$  and  $r$ , respectively, mapping the support of  $(Z, X)$  and  $X$  to  $(\epsilon, 1 - \epsilon)$  for some  $\epsilon \in (0, 1/2)$ .

### *Second key input: High-quality machine learning methods*

The second key input is the use of high-quality machine learning estimators for the nuisance parameters.

For instance, in the PLR model with `score = "IV-type"`, we need to have access to consistent estimators of  $g_0$  and  $m_0$  with respect to the  $L^2(P)$  norm  $\|\cdot\|_{P,2}$ , such that

$$\|\hat{m}_0 - m_0\|_{P,2} + \|\hat{\ell}_0 - \ell_0\|_{P,2} \leq o(N^{-1/4}).$$

In the PLIV model, the sufficient condition is

$$\|\hat{r}_0 - r_0\|_{P,2} + \|\hat{m}_0 - m_0\|_{P,2} + \|\hat{\ell}_0 - \ell_0\|_{P,2} \leq o(N^{-1/4}).$$

These conditions are plausible for many ML methods. Different structured assumptions on  $\eta_0$  lead to the use of different machine-learning tools for estimating  $\eta_0$  as listed in [Chernozhukov et al. \(2018, pp. 22–23\)](#):

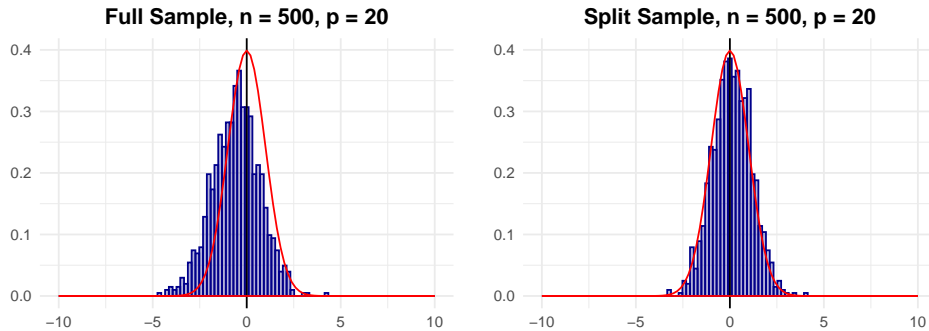


Figure 4: Performance of orthogonal estimators based on full sample and sample splitting in a simulated data example. Left panel: Histogram of the studentized estimator  $\hat{\theta}_0^{\text{nosplit}}$ .  $\hat{\theta}_0^{\text{nosplit}}$  is based on estimation of  $g_0$  and  $m_0$  with random forests and a procedure without sample-splitting: The entire data set is used for learning the nuisance terms and estimation of the orthogonal score. Data sets are simulated according to the data generating process in Section 2.2. Data generation and estimation are repeated 1000 times. Right panel: Histogram of the studentized DML estimator  $\tilde{\theta}_0$ .  $\tilde{\theta}_0$  is based on estimation of  $g_0$  and  $m_0$  with random forests and the cross-fitting described in Algorithm 2. Note that the simulated data sets and parameters of the random forest learners are identical to those underlying the left panel.

1. The assumption of approximate or exact sparsity for  $\eta_0$  with respect to some set of regressors, known as dictionary in computer science, calls for the use of sparsity-based machine learning methods, for example the lasso estimator, post-lasso,  $l_2$ -boosting, or forward selection, among others.
2. The assumption of density of  $\eta_0$  with respect to some dictionary calls for density-based estimators such as the ridge. Mixed structures based on sparsity and density suggest the use of elastic net or lava.
3. If  $\eta_0$  can be well approximated by tree-based methods, regression trees and random forests are suitable.
4. If  $\eta_0$  can be well approximated by sparse, shallow or deep neural networks,  $l_1$ -penalized neural networks, shallow neural networks or deep neural networks are attractive.

For most of these ML methods, performance guarantees are available that make it possible to satisfy the theoretical requirements. For deep learning results can be found in Farrell, Liang, and Misra (2021), for lasso in Bühlmann and Van de Geer (2011). Moreover, if  $\eta_0$  can be well approximated by at least one model mentioned in the list above, ensemble or aggregated methods (Wolpert 1992; Breiman 1996) can be used. Ensemble and aggregation methods ensure that the performance guarantee is approximately no worse than the performance of the best method (Van der Laan, Polley, and Hubbard 2007; Dudoit and Van der Laan 2005).

### *Third key input: Sample splitting*

The third key input is to use a form of sample splitting at the stage of producing the estimator of the main parameter  $\theta_0$ , which allows to avoid biases arising from overfitting.

Biases arising from overfitting could result from using highly complex fitting methods such as boosting, random forests, ensemble, and hybrid machine learning methods. We specifically use cross-fitted forms of the empirical moments, as detailed below in Algorithms 1 and 2, in

estimation of  $\theta_0$ . If the same samples would be used to estimate  $\eta_0$  and the causal parameter  $\theta_0$ , we may end up with very large bias, which we refer to as an overfitting bias. While sample splitting is key for the DML approach, other approaches, like target maximum likelihood, allow for the use of arbitrary machine learning methods for the estimation of the nuisance parameters without sample splitting. The overfitting bias is illustrated in Figure 4. The left panel shows the histogram of a studentized estimator  $\hat{\theta}_0^{\text{nosplit}}$  with  $\hat{\theta}_0^{\text{nosplit}}$  being obtained from solving the orthogonal score of Equation 11 without sample splitting. All observations are used to learn functions  $g_0$  and  $m_0$  in the PLR model and to solve the score  $\frac{1}{N} \sum_i^N \psi(W_i; \hat{\theta}_0^{\text{nosplit}}, \hat{\eta}_0)$ . Consequently, this overfitting bias leads to a considerable shift of the empirical distribution to the left. The double machine learning estimator underlying the histogram in the right panel is obtained with cross-fitting according to Algorithm 2. The sample-splitting procedure makes it possible to completely eliminate the bias induced by overfitting.

## 5. The double machine learning inference method

### 5.1. Double machine learning for estimation of a causal parameter

We assume that we have a sample  $(W_i)_{i=1}^N$ , modeled as i.i.d. copies of  $W = (Y, D, Z, X)$ , whose law is determined by the probability measure  $P$ . We assume that  $N$  is divisible by  $K$  in order to simplify the notation. Let  $E_N$  denote the empirical expectation

$$E_N[g(W)] := \frac{1}{N} \sum_{i=1}^N g(W_i).$$

*Algorithm 1: DML1 (generic double machine learning with cross-fitting)*

- (1) Inputs: Choose a model (PLR, PLIV, IRM, IIVM), provide data  $(W_i)_{i=1}^N$ , a Neyman-orthogonal score function  $\psi(W; \theta, \eta)$ , which depends on the model being estimated, and specify machine learning methods for  $\eta$ .
- (2) Train ML predictors on folds: Take a  $K$ -fold random partition  $(I_k)_{k=1}^K$  of observation indices  $[N] = \{1, \dots, N\}$  such that the size of each fold  $I_k$  is  $n = N/K$ . For each  $k \in [K] = \{1, \dots, K\}$ , construct a high-quality machine learning estimator

$$\hat{\eta}_{0,k} = \hat{\eta}_{0,k}((W_i)_{i \notin I_k})$$

of  $\eta_0$ , where  $x \mapsto \hat{\eta}_{0,k}(x)$  depends only on the subset of data  $(W_i)_{i \notin I_k}$ .

- (3) For each  $k \in [K]$ , construct the estimator  $\check{\theta}_{0,k}$  as the solution to the equation

$$\frac{1}{n} \sum_{i \in I_k} \psi(W_i; \check{\theta}_{0,k}, \hat{\eta}_{0,k}) = 0.$$

The estimate of the causal parameter is obtained via aggregation

$$\tilde{\theta}_0 = \frac{1}{K} \sum_{k=1}^K \check{\theta}_{0,k}.$$



- (4) Output: The estimate of the causal parameter  $\tilde{\theta}_0$  as well as the values of the evaluated score function are returned.

*Algorithm 2: DML2 (generic double machine learning with cross-fitting)*

- (1) Inputs: Choose a model (PLR, PLIV, IRM, IIVM), provide data  $(W_i)_{i=1}^N$ , a Neyman-orthogonal score function  $\psi(W; \theta, \eta)$ , which depends on the model being estimated, and specify machine learning methods for  $\eta$ .
- (2) Train ML predictors on folds: Take a  $K$ -fold random partition  $(I_k)_{k=1}^K$  of observation indices  $[N] = \{1, \dots, N\}$  such that the size of each fold  $I_k$  is  $n = N/K$ . For each  $k \in [K] = \{1, \dots, K\}$ , construct a high-quality machine learning estimator

$$\hat{\eta}_{0,k} = \hat{\eta}_{0,k}((W_i)_{i \notin I_k})$$

of  $\eta_0$ , where  $x \mapsto \hat{\eta}_{0,k}(x)$  depends only on the subset of data  $(W_i)_{i \notin I_k}$ .

- (3) Construct the estimator for the causal parameter  $\tilde{\theta}_0$  as the solution to the equation

$$\frac{1}{N} \sum_{k=1}^K \sum_{i \in I_k} \psi(W_i; \tilde{\theta}_0, \hat{\eta}_{0,k}) = 0.$$

- (4) Output: The estimate of the causal parameter  $\tilde{\theta}_0$  as well as the values of the evaluated score function are returned.

Both Algorithm 1 and 2 use out-of-sample predictions generated by ML learners in order to solve an orthogonal moment condition and, hence, share the same steps (1) and (2). However, the algorithms differ in the way the nuisance predictions are plugged into the score function and in the subsequent solution for  $\theta_0$ . In Algorithm 1, the score is solved on each of the  $K$  folds and the estimate  $\tilde{\theta}_0$  is obtained by averaging the  $K$  preliminary estimators,  $\check{\theta}_{0,k}$  with  $k = 1, \dots, K$ . According to Algorithm 2, the out-of-sample predictions  $\hat{\eta}_0$  are all plugged into one score function, which is then solved to obtain the estimate  $\tilde{\theta}_0$ .

*Remark 1: Linear scores*

The score for the models PLR, PLIV, IRM and IIVM are linear in  $\theta$ , having the form

$$\psi(W; \theta, \eta) = \psi_a(W; \eta)\theta + \psi_b(W; \eta),$$

hence the estimator  $\tilde{\theta}_{0,k}$  for DML2 ( $\check{\theta}_{0,k}$  for DML1) takes the form

$$\tilde{\theta}_0 = -(\mathbf{E}_N[\psi_a(W; \eta)])^{-1} \mathbf{E}_N[\psi_b(W; \eta)].$$

The linear score function representations of the PLR, PLIV, IRM and IIVM are  
 PLR with score = "partialling out"

$$\begin{aligned} \psi_a(W; \eta) &= -(D - m(X))(D - m(X)), \\ \psi_b(W; \eta) &= (Y - \ell(X))(D - m(X)). \end{aligned}$$

PLR with score = "IV-type"

$$\begin{aligned}\psi_a(W; \eta) &= -D(D - m(X)), \\ \psi_b(W; \eta) &= (Y - g(X))(D - m(X)).\end{aligned}$$

PLIV with score = "partialling out"

$$\begin{aligned}\psi_a(W; \eta) &= -(D - r(X))(Z - m(X)), \\ \psi_b(W; \eta) &= (Y - \ell(X))(Z - m(X)).\end{aligned}$$

PLIV with score = "IV-type"

$$\begin{aligned}\psi_a(W; \eta) &= -D(Z - m(X)), \\ \psi_b(W; \eta) &= (Y - g(X))(Z - m(X)).\end{aligned}$$

IRM with score = "ATE"

$$\begin{aligned}\psi_a(W; \eta) &= -1, \\ \psi_b(W; \eta) &= g(1, X) - g(0, X) + \frac{D(Y - g(1, X))}{m(X)} - \frac{(1 - D)(Y - g(0, X))}{1 - m(x)}.\end{aligned}$$

IRM with score = "ATTE"

$$\begin{aligned}\psi_a(W; \theta, \eta) &= -\frac{D}{p} \\ \psi_b(W; \theta, \eta) &= \frac{D(Y - g(0, X))}{p} - \frac{m(X)(1 - D)(Y - g(0, X))}{p(1 - m(x))}\end{aligned}$$

IIVM with score = "LATE"

$$\begin{aligned}\psi_a(W; \eta) &= -\left(r(1, X) - r(0, X) + \frac{Z(D - r(1, X))}{m(X)} - \frac{(1 - Z)(D - r(0, X))}{1 - m(x)}\right), \\ \psi_b(W; \eta) &= g(1, X) - g(0, X) + \frac{Z(Y - g(1, X))}{m(X)} - \frac{(1 - Z)(Y - g(0, X))}{1 - m(x)}.\end{aligned}$$

*Remark 2: Sample splitting*

In Step (2) of the Algorithm DML1 and DML2, the estimator  $\hat{\eta}_{0,k}$  can generally be an ensemble or aggregation of several estimators as long as we only use the data  $(W_i)_{i \notin I_k}$  outside the  $k$ -th fold to construct the estimators.

*Remark 3: Recommendation*

We have found that  $K = 4$  or  $K = 5$  to work better than  $K = 2$  in a variety of empirical examples and in simulations. The default for the option `n_folds` that implements the value of  $K$  is `n_folds=5`. Moreover, we generally recommend to repeat the estimation procedure multiple times and use the estimates and standard errors as aggregated over multiple repetitions as described in [Chernozhukov et al. \(2018, C30-C31\)](#). This aggregation will be automatically executed if the number of repetitions `n_rep` is set to a value larger than 1.

The properties of the estimator are as follows.

**Theorem 1** *There exist regularity conditions, such that the estimator  $\tilde{\theta}_0$  concentrates in a  $1/\sqrt{N}$ -neighborhood of  $\theta_0$  and the sampling error  $\sqrt{N}(\tilde{\theta}_0 - \theta_0)$  is approximately normal*

$$\sqrt{N}(\tilde{\theta}_0 - \theta_0) \rightsquigarrow N(0, \sigma^2),$$

with mean zero and variance given by

$$\begin{aligned}\sigma^2 &= J_0^{-2} \mathbf{E}(\psi^2(W; \theta_0, \eta_0)), \\ J_0 &= \mathbf{E}(\psi_a(W; \eta_0)).\end{aligned}$$

*Algorithm 3: Variance estimation and confidence intervals*

- (1) Inputs: Use the inputs and outputs from Algorithm 1 (DML1) or Algorithm 2 (DML2).
- (2) Variance and confidence intervals: Estimate the asymptotic variance of  $\tilde{\theta}_0$  by

$$\begin{aligned}\hat{\sigma}^2 &= \hat{J}_0^{-2} \frac{1}{N} \sum_{k=1}^K \sum_{i \in I_k} [\psi(W_i; \tilde{\theta}_0, \hat{\eta}_{0,k})]^2, \\ \hat{J}_0 &= \frac{1}{N} \sum_{k=1}^K \sum_{i \in I_k} \psi_a(W_i; \hat{\eta}_{0,k})\end{aligned}$$

and form an approximate  $(1 - \alpha)$  confidence interval, which is asymptotically valid, as

$$[\tilde{\theta}_0 \pm \Phi^{-1}(1 - \alpha/2) \hat{\sigma} / \sqrt{N}].$$

- (3) Output: Output variance estimator and the confidence interval.

**Theorem 2** *Under the same regularity condition, this interval contains  $\theta_0$  for approximately  $(1 - \alpha) \times 100$  percent of data realizations*

$$\mathbf{P} \left( \theta_0 \in \left[ \tilde{\theta}_0 \pm \Phi^{-1}(1 - \alpha/2) \hat{\sigma} / \sqrt{N} \right] \right) \rightarrow (1 - \alpha).$$

*Remark 4: Brief literature overview on double machine learning*

The presented double machine learning method was developed in Chernozhukov *et al.* (2018). The idea of using property 10 to construct estimators and inference procedures that are robust to small mistakes in nuisance parameters can be traced back to Neyman (1959) and has been used explicitly or implicitly in the literature on debiased sparsity-based inference (Belloni, Chernozhukov, and Hansen 2011; Belloni, Chernozhukov, and Wang 2014b; Javanmard and Montanari 2014; Van de Geer, Bühlmann, Ritov, and Dezeure 2014; Zhang and Zhang 2014; Chernozhukov, Hansen, and Spindler 2015b) as well as (implicitly) in the classical semi-parametric learning theory with low-dimensional  $X$  (Levit 1975; Hasminskii and Ibragimov 1978; Bickel, Klaassen, Ritov, and Wellner 1993; Newey 1994; Van der Vaart 2000; Van der Laan and Rose 2011). These references also explain that if we use scores  $\psi$  that are not Neyman-orthogonal in high dimensional settings, then the resulting estimators of  $\theta_0$  are not  $1/\sqrt{N}$  consistent and are generally heavily biased.

*Remark 5: Literature on sample splitting*

Sample splitting has been used in the traditional semiparametric estimation literature to establish good properties of semiparametric estimators under weak conditions (Klaassen 1987; Schick 1986; Van der Vaart 2000; Zheng and Van der Laan 2011). In sparse learning problems with high-dimensional  $X$ , sample splitting was employed in Belloni, Chen, Chernozhukov, and Hansen (2012). There and here, the use of sample splitting results in weak conditions on the estimators of nuisance parameters, translating into weak assumptions on sparsity in the case of sparsity-based learning.

*Remark 6: Debiased machine learning*

The presented approach builds upon and generalizes the approach of Belloni *et al.* (2011), Zhang and Zhang (2014), Javanmard and Montanari (2014), Javanmard and Montanari (2014), Javanmard and Montanari (2018), Belloni, Chernozhukov, and Hansen (2014c), Belloni, Chernozhukov, and Kato (2014a), Bühlmann and Van de Geer (2015), which considered estimation of the special case (1)–(2) using lasso without cross-fitting. This generalization, by relying upon cross-fitting, opens up the use of a much broader collection of machine learning methods and, in the case the lasso is used to estimate the nuisance functions, allows relaxation of sparsity conditions. All of these approaches can be seen as “debiasing” the estimation of the main parameter by constructing, implicitly or explicitly, score functions that satisfy the exact or approximate Neyman orthogonality.

**5.2. Methods for simultaneous inference**

In addition to estimation of target causal parameters, standard errors, and confidence intervals, the package **DoubleML** provides methods to perform valid simultaneous inference based on a multiplier bootstrap procedure introduced in Chernozhukov, Chetverikov, and Kato (2013) and Chernozhukov, Chetverikov, and Kato (2014) and suggested in high-dimensional linear regression models in Belloni *et al.* (2014a). Accordingly, it is possible to (i) construct simultaneous confidence bands for a potentially large number of causal parameters and (ii) adjust  $p$  values in a test of multiple hypotheses based on the inferential procedure introduced above.

We consider a causal PLR with  $p_1$  causal parameters of interest  $\theta_{0,1}, \dots, \theta_{0,p_1}$  associated with the treatment variables  $D_1, \dots, D_{p_1}$ . The parameter of interest  $\theta_{0,j}$  with  $j = 1, \dots, p_1$  solves a corresponding moment condition

$$\mathbb{E}[\psi_j(W; \theta_{0,j}, \eta_{0,j})] = 0,$$

as for example considered in Belloni, Chernozhukov, Chetverikov, and Wei (2018). To perform inference in a setting with multiple target coefficients  $\theta_{0,j}$ , the double machine learning procedure implemented in **DoubleML** iterates over the target variables of interest. During estimation of the effect of treatment  $D_j$  on  $Y$  as measured by the coefficient  $\theta_{0,j}$ , the remaining treatment variables enter the nuisance terms by default, i.e., they are added to the set of control variables  $X$ .

*Algorithm 4: Multiplier bootstrap*

- (1) Inputs: Use the inputs and outputs from Algorithm 1 (DML1) or Algorithm 2 (DML2) and Algorithm 3 (Variance estimation) resulting in estimates  $\tilde{\theta}_{0,1}, \dots, \tilde{\theta}_{0,p_1}$ , and standard errors  $\hat{\sigma}_1, \dots, \hat{\sigma}_{p_1}$ .
- (2) Multiplier bootstrap: Generate random weights  $\xi_i^b$  for each bootstrap repetition  $b = 1, \dots, B$  according to a normal (Gaussian) bootstrap, wild bootstrap or exponential bootstrap. Based on the estimated standard errors given by  $\hat{\sigma}_j$  and  $\hat{J}_{0,j} = \mathbf{E}_N(\psi_{a,j}(W; \eta_{0,j}))$ , we obtain bootstrapped versions of the  $t$  statistics  $t_j^{*,b}$  for  $j = 1, \dots, p_1$

$$t_j^{*,b} = \frac{1}{\sqrt{N} \hat{J}_{0,j} \hat{\sigma}_j} \sum_{k=1}^K \sum_{i \in I_k} \xi_i^b \cdot \psi_j(W_i; \tilde{\theta}_{0,j}, \hat{\eta}_{0,j;k}).$$

- (3) Output: Output the bootstrapped test statistics.

*Remark 7: Computational efficiency*

The multiplier bootstrap procedure of Chernozhukov *et al.* (2013) and Chernozhukov *et al.* (2014) is computationally efficient because it does not require resampling and reestimation of the causal parameters. Instead, it is sufficient to introduce a random perturbation of the score  $\psi$  and solve for  $\theta_0$ , accordingly.

To construct simultaneous  $(1 - \alpha)$ -confidence bands, the multiplier bootstrap presented in Algorithm 4 can be used to obtain a constant  $c_{1-\alpha}$  that will guarantee asymptotic  $(1 - \alpha)$  coverage

$$\left[ \tilde{\theta}_{0,j} \pm c_{1-\alpha} \cdot \hat{\sigma}_j / \sqrt{N} \right]. \quad (12)$$

The constant  $c_{1-\alpha}$  is obtained in two steps.

1. Calculate the maximum of the absolute values of the bootstrapped  $t$  statistics,  $t_j^{*,b}$ , in every repetition  $b$  with  $b = 1, \dots, B$ .
2. Use the  $(1 - \alpha)$ -quantile of the  $B$  maxima statistics from Step 1 as  $c_{1-\alpha}$  and construct simultaneous confidence bands according to Equation 12.

Moreover, it is possible to derive an adjustment method for  $p$  values obtained from a test of multiple hypotheses, including classical adjustments such as the Bonferroni correction as well as the Romano-Wolf stepdown procedure (Romano and Wolf 2005a,b). The latter is implemented according to the algorithm for adjustment of  $p$  values as provided in Romano and Wolf (2016) and adapted to high-dimensional linear regression based on the lasso in Bach, Chernozhukov, and Spindler (2018).

## 6. Implementation details

In this section, we briefly provide information on the implementation details such as the class structure, the data backend and the use of machine learning methods. Section 7 provides a

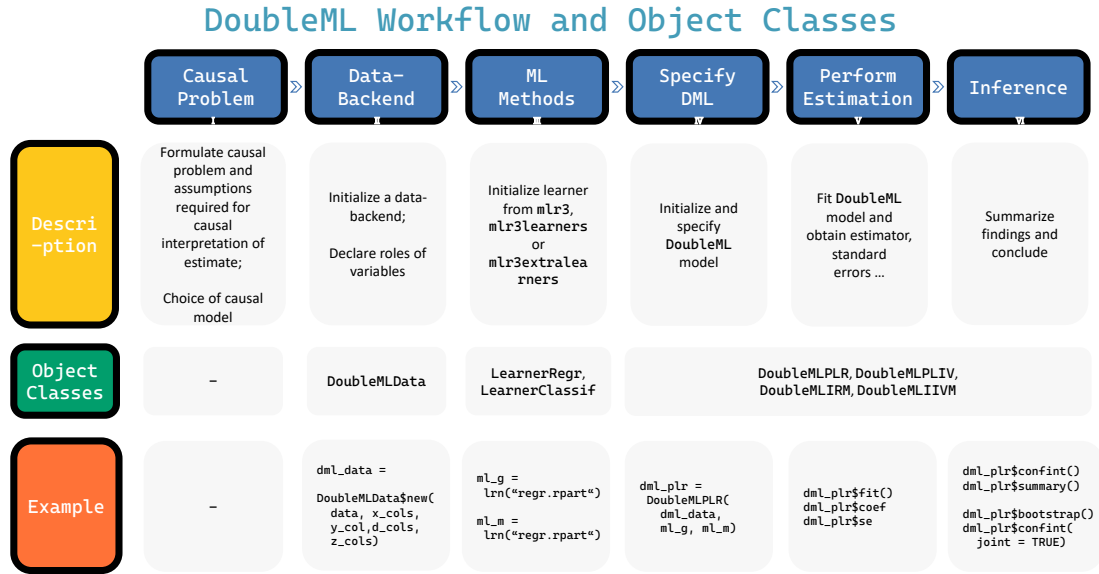
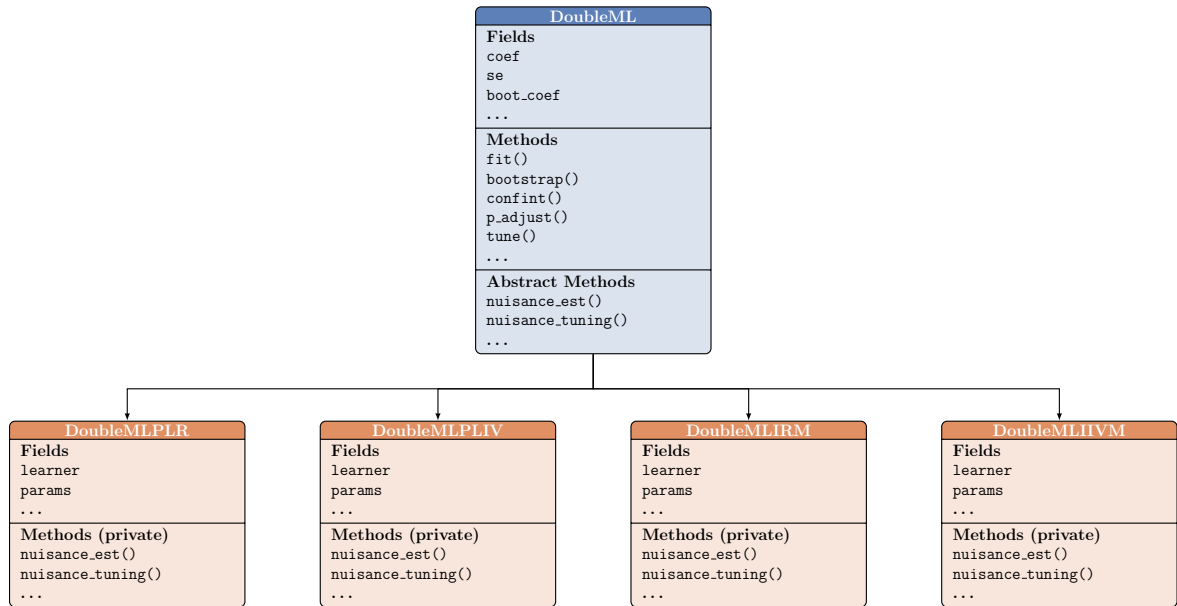


Figure 5: Flowchart with main steps and object classes in **DoubleML**. The flowchart illustrates the basic steps for estimation of causal parameters with **DoubleML**. The diagram contains a short description of the main steps and lists the object classes used in each step. A short example demonstrates the use of the object classes and methods.

demonstration of **DoubleML** in real-data and simulation examples. More information on the implementation can be found in the DoubleML User Guide, that is available online at <https://docs.doubleml.org/stable/>. All class methods are documented in the documentation of the corresponding class, which can be browsed online at <https://docs.doubleml.org/r/stable/> or, for example, by using the commands `help(DoubleML)`, `help(DoubleMLPLR)`, or `help(DoubleMLData)` in R. For an introduction to **R6** we refer to the introduction of the online book for **mlr3**, available at <https://mlr3book.mlr-org.com/intro.html>.

## 6.1. Object orientation and class structure

As pointed out in the previous sections, the double machine learning framework provides a general inferential framework in that it covers a plethora of causal models that can be characterized in terms of a Neyman-orthogonal score function  $\psi$ . In order to design an implementation that is similarly general, the implementation of **DoubleML** for R is based on object orientation as enabled by the the **R6** package (Chang 2021). The choice of the object orientation provided by **R6** as compared to alternative approaches (e.g., **S3** or **S4** classes) has been motivated by mainly three reasons: First, we would like to obtain an optimal compatibility with the **mlr3** ecosystem that is built with **R6** classes as well. Second, **R6** makes it possible to use encapsulation, inheritance, active bindings and to distinguish between private and public methods which are important features required in our implementation. Third, the object-oriented implementation of **DoubleML** makes it possible to achieve a high degree of comparability with its Python twin, which will likely facilitate and accelerate the continuous development of both packages in the future. For an introduction to object orientation in R and the **R6** package, we refer to the vignettes of the **R6** package that are available online

Figure 6: Class structure of the **DoubleML** package for R.

at <https://r6.r-lib.org/articles/>, Chapter 2.1 of Becker, Binder, Bischl, Lang, Pfisterer, Reich, Richter, Schratz, and Sonabend (2021), and the chapters on object orientation in Wickham (2019). The structure of the classes are presented in Figure 6. Moreover, the flowchart in Figure 5 illustrates the main steps of an analysis in **DoubleML** and links them to the provided object classes. Figure 5 provides a short code demonstration, too. The abstract class ‘**DoubleML**’ provides all methods for estimation and inference, for example the methods `fit()`, `bootstrap()`, `confint()`. All key components associated with estimation and inference are implemented in ‘**DoubleML**’, for example the sample splitting, the implementation of Algorithm 1 (DML1) and Algorithm 2 (DML2), the estimation of the causal parameters, and the computation of the scores  $\psi(W; \theta, \eta)$ . Only the model-specific properties and methods are allocated at the classes ‘**DoubleMLPLR**’ (implementing the PLR), ‘**DoubleMLPLIV**’ (PLIV), ‘**DoubleMLIRM**’ (IRM), and ‘**DoubleMLIIVM**’ (IIVM). For example, each of the models has one or several Neyman-orthogonal score functions that are implemented for the specific child classes.

## 6.2. Data backend and causal model

The ‘**DoubleMLData**’ class serves as the data backend and implements the causal model of interest. The user is required to specify the roles of the variables in a data set at hand. Depending on the causal model considered, it is necessary to declare the dependent variable, the treatment variable(s), confounding variables(s), and, in the case of instrumental variable regression, one or multiple instruments. The data backend can be initialized from a `data.table` (Dowle and Srinivasan 2023). **DoubleML** provides wrappers to initialize from ‘`data.frame`’ and ‘`matrix`’ objects, as well.

### 6.3. Learners, parameters and tuning

Generally, all learners provided by the packages **mlr3**, **mlr3learners** and **mlr3extralearners** can be used for estimation of the nuisance functions of the structural models presented above. An interactive list of supported learners is available at the **mlr3extralearners** website ([https://mlr3extralearners.ml-org.com/articles/learners/list\\_learners.html](https://mlr3extralearners.ml-org.com/articles/learners/list_learners.html)). The **mlr3extralearners** package makes it possible to add new learners, as well. The performance of the double machine learning estimator  $\tilde{\theta}_0$  will depend on the predictive quality of the used machine learning method. Machine learning methods usually have several (hyper-)parameter that need to be adapted to the specific application. Tuning of model parameters can be either performed externally or internally. The latter is implemented in the method `tune()` and is further illustrated in an example in Section 7.6. Both cases build on the functionalities provided by the package **mlr3tuning**.

### 6.4. Modifications and extensions

The flexible architecture of the **DoubleML** package allows users to modify the estimation procedure in many regards. Among others, users can provide customized sample splitting rules after initialization of the causal model via the method `set_sample_splitting()`. Moreover, it is possible to adjust the Neyman-orthogonal score function by externally providing a customized function via the `score` option during initialization of the causal model object. Short examples for both of these potential extensions are presented in Section 7.7.

## 7. Estimation in real-data and simulated examples

In this section, we will first demonstrate the use of **DoubleML** in a real-data example, which is based on data from the Pennsylvania Reemployment Bonus experiment (Bilias 2000). This empirical example has been used in Chernozhukov *et al.* (2018), as well. The goal in the empirical example is to estimate the causal parameter in a partially linear and an interactive regression model. We further provide a short example on how valid simultaneous inference can be performed with **DoubleML**. Finally, we present results from a short simulation study as a brief assessment of the finite-sample performance of the implemented estimators. Here we want to stress that in real world applications modelling choices of the estimation of the nuisance parameters and proper tuning of the parameters are very important. We would like to mention that the presented examples are mainly included for the purpose of illustration. In practice, we recommend to carefully choose and tune the ML learners in terms of their hyperparameters.

### 7.1. Initialization of the data backend

We begin our real-data example by downloading the Pennsylvania Reemployment Bonus data set. To do so, we use the call `fetch_bonus()` (a connection to the internet is required).

```
R> library("DoubleML")
```

Load data as `data.table`.

```
R> dt_bonus <- fetch_bonus(return_type = "data.table")
```



The output is suppressed for the sake of brevity.

```
R> dt_bonus
```

The data backend ‘DoubleMLData’ can be initialized from a ‘data.table’ object by specifying the dependent variable  $Y$  via a character in `y_col`, the treatment variable(s)  $D$  in `d_cols`, and the confounders  $X$  via `x_cols`. Moreover, in IV models, an instrument can be specified via `z_cols`. In the next step, we assign the roles to the variables in the data set: `y_col = 'inuidur1'` serves as outcome variable  $Y$ , the column `d_cols = 'tg'` serves as treatment variable  $D$  and the columns `x_cols` specify the confounders.

```
R> obj_dml_data_bonus <- DoubleMLData$new(dt_bonus,
+   y_col = "inuidur1",
+   d_cols = "tg",
+   x_cols = c("female", "black", "othrace", "dep1", "dep2", "q2", "q3",
+   "q4", "q5", "q6", "age<35", "age>54", "durable", "lusc", "husd"))
```

The data backend can be printed to obtain a summary of the main attributes of the ‘DoubleMLData’ object.

```
R> obj_dml_data_bonus
```

```
===== DoubleMLData Object =====
----- Data summary -----
Outcome variable: inuidur1
Treatment variable(s): tg
Covariates: female, black, othrace, dep1, dep2, q2, q3, q4, q5, q6, age<35,
  age>54, durable, lusc, husd
Instrument(s):
No. Observations: 5099
```

Print the data set (output suppressed).

```
R> obj_dml_data_bonus$data
```

#### *Remark 8: Wrappers for the data backend*

To initialize an instance of the class ‘DoubleMLData’ from a ‘data.frame’ or a collection of ‘matrix’ objects, **DoubleML** provides the convenient wrappers `double_ml_data_from_data_frame()` and `double_ml_data_from_matrix()`. Although the data backend does not provide a formula interface, ‘DoubleMLData’ objects can be initialized from a ‘model.matrix’ object. The following example demonstrates how users may proceed to specify the causal model by using a ‘formula’. We load the Pennsylvania Reemployment Bonus data set as a ‘data.frame’ and replicate a flexible model specification used in the empirical analysis of [Chernozhukov et al. \(2018, C38–C40\)](#). To flexibly model the nuisance function, we generate all two-way interactions of the control variables.

Load the data as a ‘data.frame’.

```
R> df_bonus <- fetch_bonus(return_type = "data.frame")
```

Print the names of the variables.

```
R> names(df_bonus)
```

```
[1] "inuidur1" "female" "black" "othrace" "dep1" "dep2"
[7] "q2"      "q3"      "q4"      "q5"      "q6"      "agelt35"
[13] "agegt54" "durable" "lUSD"   "hUSD"    "tg"
```

Specify a ‘formula’ object to generate all two-way interactions of the control variables.

```
R> f_flex <- formula(" ~ -1 + (female + black + othrace + dep1 + q2 + q3 +
+   q4 + q5 + q6 + agelt35 + agegt54 + durable + lUSD + hUSD)^2")
```

Create a ‘model.matrix’ based on the ‘formula’ object.

```
R> X_flex <- model.matrix(f_flex, data = df_bonus)
```

Initialize using the wrapper `double_ml_data_from_data_frame()`.

```
R> df_bonus_flex <- data.frame("inuidur1" = df_bonus$inuidur1, X_flex,
+   "tg" = df_bonus$tg)
R> obj_dml_data_bonus_flex <- double_ml_data_from_data_frame(df_bonus_flex,
+   y_col = "inuidur1", d_cols = "tg")
```

Alternatively, initialize via the wrapper `double_ml_data_from_matrix()`.

```
R> obj_dml_data_bonus_flex2 <- double_ml_data_from_matrix(X = X_flex,
+   y = df_bonus$inuidur1, d = df_bonus$tg)
```

## 7.2. Initialization of the causal model

To initialize a PLR model, we have to provide a learner for each nuisance part in the model in Equation 1–2. In R, this is done by providing learners to the arguments `ml_m` for nuisance part  $m$  and `ml_l` for nuisance part  $l$ . We can pass a learner as instantiated in **mlr3** and **mlr3learners**, for example a random forest as provided by the R package **ranger** (Wright and Ziegler 2017). Previous installation of **ranger** is required. Moreover, we can specify the score (allowed choices for PLR are "partialling out" or "IV-type") and the algorithm via the option `dml_procedure` (allowed choices "dml1" and "dml2"). Optionally, it is possible to change the number of folds used for sample splitting through `n_folds` and the number of repetitions via `n_rep`, if the sample splitting and estimation procedure should be repeated.

Set a seed for replication of the sample split.

```
R> set.seed(31415)
R> learner_l <- lrn("regr.ranger", num.trees = 500, min.node.size = 2,
+   max.depth = 5)
```

```

R> learner_m <- lrn("regr.ranger", num.trees = 500, min.node.size = 2,
+   max.depth = 5)
R> doubleml_bonus <- DoubleMLPLR$new(obj_dml_data_bonus, ml_l = learner_l,
+   ml_m = learner_m, score = "partialling out", dml_procedure = "dml1",
+   n_folds = 5, n_rep = 1)
R> doubleml_bonus

===== DoubleMLPLR Object =====

----- Data summary -----
Outcome variable: inuidur1
Treatment variable(s): tg
Covariates: female, black, othrace, dep1, dep2, q2, q3, q4, q5, q6, age1t35,
  age1t54, durable, lused, husd
Instrument(s):
No. Observations: 5099

----- Score & algorithm -----
Score function: partialling out
DML algorithm: dml1

----- Machine learner -----
ml_l: regr.ranger
ml_m: regr.ranger

----- Resampling -----
No. folds: 5
No. repeated sample splits: 1
Apply cross-fitting: TRUE

----- Fit summary -----

```

### 7.3. Estimation of the causal parameter in a PLR model

To perform estimation, call the `fit()` method. The output can be summarized using the `method summary()`.

```

R> doubleml_bonus$fit()
R> doubleml_bonus$summary()

```

```

Estimates and significance testing of the effect of target variables
  Estimate. Std. Error t value Pr(>|t|)
tg -0.07438    0.03543  -2.099   0.0358 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Hence, there is evidence to reject the null hypothesis that  $\theta_{0,tg} = 0$  at the 5% significance level. The estimated coefficient and standard errors can be accessed via the attributes `coef` and `se` of the object `doubleml_bonus`.

```
R> doubleml_bonus$coef
```

```
      tg
-0.07438411
```

```
R> doubleml_bonus$se
```

```
      tg
0.03543316
```

After completed estimation, we can access the resulting score  $\psi(W_i; \tilde{\theta}_0, \hat{\eta}_0)$  or the components  $\psi_a(W_i; \hat{\eta}_0)$  and  $\psi_b(W_i; \hat{\eta}_0)$ . The estimated score for the first 5 observations can be obtained via the public field `psi`. `psi` is an array with `dim = c(n_obs, n_rep, n_treat)` with

- `n_obs`: number of observations in the data,
- `n_rep`: number of repetitions (sample splitting),
- `n_treat`: number of treatment variables.

```
R> doubleml_bonus$psi[1:5, 1, 1]
```

```
[1] -0.2739454  0.7444154 -0.4509358  0.1813111 -0.3699474
```

Similarly, the components of the score  $\psi_a(W_i; \hat{\eta}_0)$  and  $\psi_b(W_i; \hat{\eta}_0)$  are available as fields.

```
R> doubleml_bonus$psi_a[1:5, 1, 1]
```

```
[1] -0.0981220 -0.1353987 -0.1276526 -0.4272341 -0.1126174
```

```
R> doubleml_bonus$psi_b[1:5, 1, 1]
```

```
[1] -0.2812441  0.7343439 -0.4604311  0.1495317 -0.3783243
```

To construct a  $(1 - \alpha)$  confidence interval, we use the `confint()` method.

```
R> doubleml_bonus$confint(level = 0.95)
```

```
      2.5 %      97.5 %
tg -0.1438318 -0.004936395
```

#### 7.4. Estimation of the causal parameter in an IRM model

The treatment variable  $D$  in the Pennsylvania Reemployment Bonus example is binary. Accordingly, it is possible to estimate an IRM model. Since the IRM requires estimation of the propensity score  $P(D | X)$ , we have to specify a classifier for the nuisance part  $m_0$ .

Initialize a classifier for estimation of the propensity score and create a new instance of a causal model, here an IRM.

```
R> learner_g <- lrn("regr.ranger", num.trees = 500, min.node.size = 2,
+   max.depth = 5)
R> learner_classif_m <- lrn("classif.ranger", num.trees = 500,
+   min.node.size = 2, max.depth = 5)
R> doubleml_irm_bonus <- DoubleMLIRM$new(obj_dml_data_bonus,
+   ml_g = learner_g, ml_m = learner_classif_m, score = "ATE",
+   dml_procedure = "dml1", n_folds = 5, n_rep = 1)
```

The output is suppressed for the sake of brevity.

```
R> doubleml_irm_bonus
```

To perform estimation, call the `fit()` method. The output can be summarized using the method `summary()`.

```
R> doubleml_irm_bonus$fit()
R> doubleml_irm_bonus$summary()
```

Estimates and significance testing of the effect of target variables

```
Estimate. Std. Error t value Pr(>|t|)
tg -0.07193 0.03554 -2.024 0.043 *
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The estimated coefficient is very similar to the estimate of the PLR model and our conclusions remain unchanged.

## 7.5. Simultaneous inference in a simulated data example

We consider a simulated example of a PLR model to illustrate the use of methods for simultaneous inference. First, we will generate a sparse linear model with only three variables having a non-zero effect on the dependent variable.

```
R> set.seed(3141)
R> n_obs <- 500
R> n_vars <- 100
R> theta <- rep(3, 3)
```

Generate a `data.frame` and use the corresponding wrapper.

```
R> X <- matrix(stats::rnorm(n_obs * n_vars), nrow = n_obs, ncol = n_vars)
R> y <- X[, 1:3, drop = FALSE] %*% theta + stats::rnorm(n_obs)
R> df <- data.frame(y, X)
```

We use the wrapper `double_ml_data_from_data_frame()` to specify a data backend that assigns the first 10 columns of  $X$  as treatment variables and declares the remaining columns as confounders.

```
R> doubleml_data <- double_ml_data_from_data_frame(df, y_col = "y",
+   d_cols = c("X1", "X2", "X3", "X4", "X5", "X6", "X7", "X8", "X9", "X10"))
```

Set treatment variable `d` to `X1`.

The output is suppressed for the sake of brevity.

```
R> doubleml_data
```

A sparse setting suggests the use of the lasso learner. Here, we use the lasso estimator with cross-validated choice of the penalty parameter  $\lambda$  as provided in the `glmnet` package for R (Friedman *et al.* 2010).

Output messages during fitting are suppressed.

```
R> ml_l <- lrn("regr.cv_glmnet", s = "lambda.min")
R> ml_m <- lrn("regr.cv_glmnet", s = "lambda.min")
R> doubleml_plr <- DoubleMLPLR$new(doubleml_data, ml_l, ml_m)
R> doubleml_plr$fit()
R> doubleml_plr$summary()
```

Estimates and significance testing of the effect of target variables

	Estimate.	Std. Error	t value	Pr(> t )	
X1	3.017802	0.046180	65.348	<2e-16	***
X2	3.025812	0.042683	70.891	<2e-16	***
X3	3.000914	0.045849	65.452	<2e-16	***
X4	-0.034815	0.040955	-0.850	0.3953	
X5	0.035118	0.048132	0.730	0.4656	
X6	0.002171	0.044622	0.049	0.9612	
X7	-0.036129	0.046798	-0.772	0.4401	
X8	0.020361	0.044048	0.462	0.6439	
X9	-0.019439	0.043180	-0.450	0.6526	
X10	0.076180	0.043682	1.744	0.0812	.

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

The multiplier bootstrap procedure can be executed using the `bootstrap()` method where the option `method` specifies the choice of the random perturbations and `n_rep_boot` the number of bootstrap repetitions.

```
R> doubleml_plr$bootstrap(method = "normal", n_rep_boot = 1000)
```

The resulting bootstrapped  $t$  statistics are available via the field `boot_t_stat`. To construct a simultaneous confidence interval, we set the option `joint = TRUE` when calling the `confint()` method.

```
R> doubleml_plr$confint(joint = TRUE)
```

	2.5 %	97.5 %
X1	2.88766757	3.14793595
X2	2.90553386	3.14609021
X3	2.87171334	3.13011430
X4	-0.15022399	0.08059423
X5	-0.10051468	0.17075155
X6	-0.12357302	0.12791441
X7	-0.16800517	0.09574654
X8	-0.10376590	0.14448792
X9	-0.14111984	0.10224143
X10	-0.04691574	0.19927524

The correction of the  $p$  values of a joint hypotheses test on the considered causal parameters is implemented in the method `p_adjust()`. By default, the adjustment procedure specified in the option `method` is the Romano-Wolf stepdown procedure.

```
R> doubleml_plr$p_adjust(method = "romano-wolf")
```

	Estimate.	pval
X1	3.017801759	0.000
X2	3.025812035	0.000
X3	3.000913821	0.000
X4	-0.034814877	0.942
X5	0.035118436	0.942
X6	0.002170694	0.961
X7	-0.036129317	0.942
X8	0.020361010	0.951
X9	-0.019439209	0.951
X10	0.076179750	0.451

Alternatively, the correction methods provided in the `stats` function `p.adjust` can be applied, for example the Bonferroni, Bonferroni-Holm, or Benjamini-Hochberg correction. For example a Bonferroni correction could be performed by specifying `method = "bonferroni"`.

```
R> doubleml_plr$p_adjust(method = "bonferroni")
```

	Estimate.	pval
X1	3.017801759	0.0000000
X2	3.025812035	0.0000000
X3	3.000913821	0.0000000
X4	-0.034814877	1.0000000
X5	0.035118436	1.0000000
X6	0.002170694	1.0000000
X7	-0.036129317	1.0000000
X8	0.020361010	1.0000000
X9	-0.019439209	1.0000000
X10	0.076179750	0.8116808

## 7.6. Learners, parameters and tuning

The performance of the final double machine learning estimator depends on the predictive performance of the underlying ML method. First, we briefly show how externally tuned parameters can be passed to the learners in **DoubleML**. Second, it is demonstrated how the parameter tuning can be done internally by **DoubleML**.

### *External tuning and parameter passing*

Section 3 of the **mlr3** book (Becker *et al.* 2021) provides a step-by-step introduction to the powerful tuning functionalities of the **mlr3tuning** package. Accordingly, it is possible to manually reconstruct the **mlr3** regression and classification problems, which are internally handled in **DoubleML**, and to perform parameter tuning accordingly. One advantage of this procedure is that it allows users to fully exploit the powerful benchmarking and tuning tools of **mlr3** and **mlr3tuning**.

Consider the sparse regression example from above. We will briefly consider a setting where we explicitly set the parameter  $\lambda$  for a **glmnet** estimator rather than using the internal cross-validated choice with `cv_glmnet`.

Suppose for simplicity, some external tuning procedure resulted in an optimal value of  $\lambda = 0.1$  for nuisance part  $m$  and  $\lambda = 0.09$  for nuisance part  $\ell$  for the first treatment variable and  $\lambda = 0.095$  and  $\lambda = 0.085$  for the second variable, respectively. After initialization of the model object, we can set the parameter values using the method `set_ml_nuisance_params()`.

```
R> ml_l <- lrn("regr.glmnet")
R> ml_m <- lrn("regr.glmnet")
R> doubleml_plr <- DoubleMLPLR$new(doubleml_data, ml_l, ml_m)
```

To set the values, we have to specify the treatment variable and the nuisance part. If no values are set, the default values are used. Note that variable names are overwritten by the wrapper for the matrix interface.

```
R> doubleml_plr$set_ml_nuisance_params("ml_m", "X1",
+   param = list("lambda" = 0.1))
R> doubleml_plr$set_ml_nuisance_params("ml_l", "X1",
+   param = list("lambda" = 0.09))
R> doubleml_plr$set_ml_nuisance_params("ml_m", "X2",
+   param = list("lambda" = 0.095))
R> doubleml_plr$set_ml_nuisance_params("ml_l", "X2",
+   param = list("lambda" = 0.085))
```

All externally specified parameters can be retrieved from the field `params`. The output is omitted for the sake of brevity.

```
R> str(doubleml_plr$params)
R> doubleml_plr$fit()
R> doubleml_plr$summary()
```

Estimates and significance testing of the effect of target variables  
 Estimate. Std. Error t value Pr(>|t|)



```

X1  3.041094  0.060030  50.660  <2e-16 ***
X2  2.993916  0.054590  54.844  <2e-16 ***
X3  2.993419  0.055144  54.283  <2e-16 ***
X4  -0.035201  0.040637  -0.866   0.386
X5   0.021541  0.047569   0.453   0.651
X6  -0.006652  0.044715  -0.149   0.882
X7  -0.039650  0.046823  -0.847   0.397
X8   0.011146  0.044037   0.253   0.800
X9  -0.021342  0.043237  -0.494   0.622
X10 0.084426  0.043641   1.935   0.053 .

```

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

### *Internal tuning and parameter passing*

An alternative to external tuning and parameter provisioning is to perform the tuning internally. The advantage of this approach is that users do not have to specify the underlying prediction problems manually. Instead, **DoubleML** uses the underlying data backend to ensure that the machine learning methods are tuned for the specific model under consideration and, hence, to possibly avoid mistakes. We initialize our structural model object with the learner. At this stage, we do not specify any parameters.

Load required packages for tuning and set logger to omit messages during tuning and fitting.

```

R> library("paradox")
R> library("mlr3tuning")
R> lgr::get_logger("mlr3")$set_threshold("warn")
R> lgr::get_logger("bbotk")$set_threshold("warn")
R> set.seed(1234)
R> ml_l <- lrn("regr.glmnet")
R> ml_m <- lrn("regr.glmnet")
R> doubleml_plr <- DoubleMLPLR$new(doubleml_data, ml_l, ml_m)

```

To perform parameter tuning, we provide a grid of values used for evaluation for each of the nuisance parameters. To set up a grid of values, we specify a named list with names corresponding to the learner names of the nuisance part (see method `learner_names()`). The elements in the list are objects of the class ‘ParamSet’ of the **paradox** package (Lang, Bischl, Richter, Sun, and Binder 2023b).

```

R> par_grids <- list(
+   "ml_l" = ParamSet$new(list(ParamDbl$new("lambda", lower = 0.05,
+     upper = 0.1))),
+   "ml_m" = ParamSet$new(list(ParamDbl$new("lambda", lower = 0.05,
+     upper = 0.1))))

```

The hyperparameter tuning is performed according to options passed through a named list `tune_settings`. The entries in the list specify options during parameter tuning with **mlr3tuning**:

- `terminator` is a `bbotk::Terminator` object passed to `mlr3tuning` that manages the budget to solve the tuning problem.
- `algorithm` is an object of class `'mlr3tuning::Tuner'` and specifies the tuning algorithm. Alternatively, `algorithm` can be a `character()` that is used as an argument in the wrapper `mlr3tuning` call `tnr(algorithm)`. The `'Tuner'` class in `mlr3tuning` supports grid search, random search, generalized simulated annealing and non-linear optimization.
- `rsmp_tune` is an object of class `'mlr3::Resampling'` that specifies the resampling method for evaluation, for example `rsmp("cv", folds = 5)` implements 5-fold cross-validation. `rsmp("holdout", ratio = 0.8)` implements an evaluation based on a hold-out sample that contains 20 percent of the observations. By default, 5-fold cross-validation is performed.
- `measure` is a named list containing the measures used for tuning of the nuisance components. The names of the entries must match the learner names (see method `learner_names()`). The entries in the list must either be objects of class `'mlr3::Measure'` or keys passed to `msr()`. If `measure` is not provided by the user, the mean squared error is used for regression models and the classification error for binary outcomes, by default.

In the next code chunk, the value of the parameter  $\lambda$  is tuned via grid search in the range 0.05 to 0.1 at a resolution of 11. The resulting grid has 11 equally spaced values ranging from a minimum value of 0.05 to a maximum value of 0.1. Type `generate_design_grid(par_grids$ml_1, resolution = 11)` to access the grid for nuisance function `ml_1`. To evaluate the predictive performance in both nuisance functions, the cross-validated mean squared error is used.

Provide tune settings.

```
R> tune_settings <- list(terminator = trm("evals", n_evals = 100),
+   algorithm = tnr("grid_search", resolution = 11),
+   rsmp_tune = rsmp("cv", folds = 5),
+   measure = list("ml_1" = msr("regr.mse"), "ml_m" = msr("regr.mse")))
```

With these parameters we can run the tuning by calling the `tune()` method for `'DoubleML'` objects.

Execution might take around 50 seconds.

```
R> doubleml_plr$tune(param_set = par_grids, tune_settings = tune_settings)
```

Output omitted for the sake of brevity, available in the appendix. Access tuning results for target variable `X1`.

```
R> doubleml_plr$tuning_res$X1
```

Access tuned parameters (output suppressed).

```
R> str(doubleml_plr$params)
```

Estimate model and call the `summary()` method.

```
R> doubleml_plr$fit()
R> doubleml_plr$summary()
```

Estimates and significance testing of the effect of target variables

	Estimate.	Std. Error	t value	Pr(> t )	
X1	3.028980	0.059701	50.736	<2e-16	***
X2	3.008650	0.054301	55.407	<2e-16	***
X3	2.960571	0.053082	55.773	<2e-16	***
X4	-0.037859	0.040976	-0.924	0.3555	
X5	0.030018	0.047880	0.627	0.5307	
X6	0.003451	0.044419	0.078	0.9381	
X7	-0.025875	0.046936	-0.551	0.5814	
X8	0.022008	0.044172	0.498	0.6183	
X9	-0.014251	0.043765	-0.326	0.7447	
X10	0.088653	0.043691	2.029	0.0424	*

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

By default, the parameter tuning is performed on the whole sample, for example in the case of  $K_{\text{tune}}$ -fold cross-validation, the entire sample is split into  $K_{\text{tune}}$  folds for evaluation of the cross-validated error. Alternatively, each of the  $K$  folds used in the cross-fitting procedure could be split up into  $K_{\text{tune}}$  subfolds that are then used for evaluation of the candidate models. As a result, the choice of the tuned parameters will be fold-specific. To perform fold-specific tuning, users can set the option `tune_on_folds = TRUE` when calling the method `tune()`.

## 7.7. Specifications and modifications of double machine learning

The flexible architecture of the **DoubleML** package allows users to modify the estimation procedure in many regards. We will shortly present two examples on how users can adjust the double machine learning framework to their needs in terms of the sample splitting procedure and the score function.

### *Sample splitting*

By default, **DoubleML** performs cross-fitting as presented in Algorithms 1 and 2. Alternatively, all implemented models allow a partition to be provided externally via the method `set_sample_splitting()`. Note that by setting `draw_sample_splitting = FALSE` one can prevent that a partition is drawn during initialization of the model object. The following calls are equivalent. In the first sample code, we use the standard interface and draw the sample-splitting with  $K = 4$  folds during initialization of the ‘DoubleMLPLR’ object.

First generate some data and initialize ML learners and a data backend.

```
R> learner <- lrn("regr.ranger", num.trees = 100, mtry = 20,
+   min.node.size = 2, max.depth = 5)
R> ml_l <- learner
R> ml_m <- learner
R> data <- make_plr_CCDDHNR2018(alpha = 0.5, n_obs = 100,
+   return_type = "data.table")
```

```
R> doubleml_data <- DoubleMLData$new(data, y_col = "y", d_cols = "d")
R> set.seed(314)
R> doubleml_plr_internal <- DoubleMLPLR$new(doubleml_data, ml_l, ml_m,
+   n_folds = 4)
R> doubleml_plr_internal$fit()
R> doubleml_plr_internal$summary()
```

Estimates and significance testing of the effect of target variables

```
Estimate. Std. Error t value Pr(>|t|)
d    0.4892    0.1024   4.776 1.79e-06 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

In the second sample code, we manually specify a sampling scheme using the `'mlr3::Resampling'` class. Alternatively, users can provide a nested list that has the following structure:

- The length of the outer list must match with the desired number of repetitions of the sample-splitting, i.e., `n_rep`.
- The inner list is a named list of length 2 specifying the `test_ids` and `train_ids`. The named entries `test_ids` and `train_ids` are lists of the same length,
  - `train_ids` is a list of length `n_folds` that specifies the indices of the observations used for model fitting in each fold,
  - `test_ids` is a list of length `n_folds` that specifies the indices of the observations used for calculation of the score in each fold.

Set up a task and cross-validation resampling scheme in **mlr3**.

```
R> doubleml_plr_external <- DoubleMLPLR$new(doubleml_data, ml_l, ml_m,
+   draw_sample_splitting = FALSE)
R> set.seed(314)
R> my_task <- Task$new("help task", "regr", data)
R> my_sampling <- rsmpl("cv", folds = 4)$instantiate(my_task)
R> train_ids <- lapply(1:4, function(x) my_sampling$train_set(x))
R> test_ids <- lapply(1:4, function(x) my_sampling$test_set(x))
R> smpls = list(list(train_ids = train_ids, test_ids = test_ids))
```

Structure of the specified sampling scheme.

```
R> str(smpls)
```

List of 1

```
$ :List of 2
..$ train_ids:List of 4
.. ..$ : int [1:75] 1 7 11 18 19 20 21 31 32 37 ...
.. ..$ : int [1:75] 10 15 16 22 26 35 38 40 41 46 ...
.. ..$ : int [1:75] 10 15 16 22 26 35 38 40 41 46 ...
.. ..$ : int [1:75] 10 15 16 22 26 35 38 40 41 46 ...
```

```

..$ test_ids :List of 4
.. ..$ : int [1:25] 10 15 16 22 26 35 38 40 41 46 ...
.. ..$ : int [1:25] 1 7 11 18 19 20 21 31 32 37 ...
.. ..$ : int [1:25] 3 5 6 8 17 24 25 28 29 34 ...
.. ..$ : int [1:25] 2 4 9 12 13 14 23 27 30 33 ...

```

Fit the model and summarize.

```

R> doubleml_plr_external$set_sample_splitting(smpls)
R> doubleml_plr_external$fit()
R> doubleml_plr_external$summary()

```

Estimates and significance testing of the effect of target variables

	Estimate.	Std. Error	t value	Pr(> t )
d	0.4892	0.1024	4.776	1.79e-06 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Setting the option `apply_cross_fitting = FALSE` at the instantiation of the causal model allows double machine learning being performed without cross-fitting. It results in randomly splitting the sample into two parts. The first half of the data is used for the estimation of the nuisance models with the machine learning methods and the second half for estimating the causal parameter, i.e., solution of the score. Note that cross-fitting performs well empirically and is recommended to remove bias induced by overfitting. Moreover, cross-fitting allows to exploit full efficiency: Every fold is used once for training the ML methods and once for estimation of the score (Chernozhukov *et al.* 2018, C6). A short example on the efficiency gains associated with cross-fitting is provided in Figure 7.

### Score function

Users may want to adjust the score function  $\psi(W; \theta_0, \eta_0)$ , for example, to adjust the DML estimators in terms of a re-weighting, e.g., to adjust for missing outcome via inverse probability of censoring weight (IPCW). An alternative to the choices provided in **DoubleML** is to pass a function via the argument `score` during initialization of the model object. The following examples are equivalent. In the first example, we use the score option `"partialling out"` for the PLR model whereas in the second case, we explicitly provide a function that implements the same score. The arguments used in the function refer to the internal objects that implement the theoretical quantities in Equation 11.

Use score `"partialling out"`.

```

R> set.seed(314)
R> doubleml_plr_partout <- DoubleMLPLR$new(doubleml_data, ml_l, ml_m,
+   score = "partialling out")
R> doubleml_plr_partout$fit()
R> doubleml_plr_partout$summary()

```

Estimates and significance testing of the effect of target variables

```

  Estimate. Std. Error t value Pr(>|t|)
d    0.5108    0.0959   5.326   1e-07 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

We define the function that implements the same score and specify the argument `score` accordingly. The function must return a named list with entries `psi_a` and `psi_b` to pass values for computation of the score.

Required input:

- `y`: dependent variable,
- `d`: treatment variable,
- `l_hat`: predicted values from regression of  $Y$  on  $X$ ,
- `m_hat`: predicted values from regression of  $D$  on  $X$ ,
- `g_hat`: predicted values from regression of  $Y - D \cdot \theta$  on  $X$ , can be ignored in this example,
- `smpis`: sample split under consideration, can be ignored in this example.

```

R> score_manual <- function(y, d, l_hat, m_hat, g_hat, smpis) {
+   resid_y = y - l_hat
+   resid_d = d - m_hat
+   psi_a = -1 * resid_d * resid_d
+   psi_b = resid_d * resid_y
+   psis = list(psi_a = psi_a, psi_b = psi_b)
+   return(psis)
+ }
R> set.seed(314)
R> doubleml_plr_manual <- DoubleMLPLR$new(doubleml_data, ml_l, ml_m,
+   score = score_manual)
R> doubleml_plr_manual$fit()
R> doubleml_plr_manual$summary()

```

Estimates and significance testing of the effect of target variables

```

  Estimate. Std. Error t value Pr(>|t|)
d    0.5108    0.0959   5.326   1e-07 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

## 7.8. A short simulation study

To illustrate the validity of the implemented double machine learning estimators, we perform a brief simulation study.

### *The role of cross-fitting*

As mentioned before the use of the cross-fitting Algorithms 1 (DML1) and 2 (DML2) makes it possible to use sample splitting and exploit full efficiency at the same time. To illustrate the

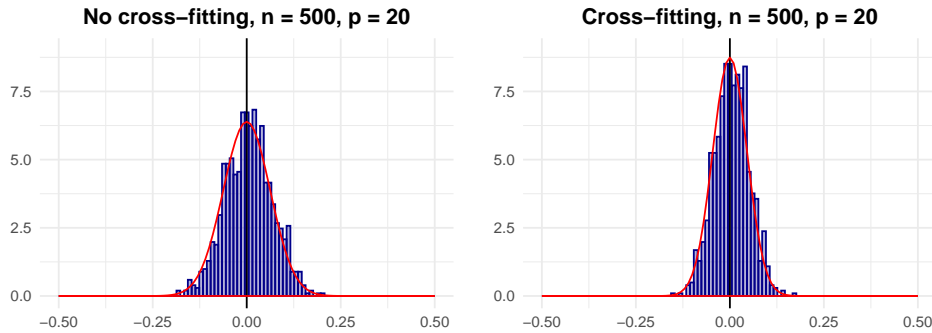


Figure 7: Illustration of efficiency gains due to the use of cross-fitting. Left panel: Histogram of the centered DML estimator without cross-fitting,  $\hat{\theta}_0^{\text{nocf}} - \theta_0$ .  $\hat{\theta}_0^{\text{nocf}}$  is the double machine learning estimator obtained from a sample split into two folds. One fold is used for estimation of the nuisance parameters and the second fold is used for evaluation of the score function and estimation. The empirical distribution can be well-approximated by a normal distribution as indicated by the red curve. Right panel: Histogram of the centered DML estimator with cross-fitting,  $\tilde{\theta}_0 - \theta_0$ . The estimator is obtained from a split into two folds and application of Algorithm 2 (DML2). In both cases, the estimators are based on estimation of  $g_0$  and  $m_0$  with random forests and an orthogonal score function provided in Equation 11. Moreover, exactly the same data sets and exactly the same partitions are used for sample splitting. The empirical distribution of the estimator that is based on cross-fitting exhibits a more pronounced concentration around zero, which reflects the smaller standard errors.

superior performance due to cross-fitting, we compare the double machine learning estimator with and without a cross-fitting procedure in the simulation setting that was presented in Section 4.1. Figure 7 illustrates that efficiency gains can be achieved if the role of the random partitions is swapped in the estimation procedure. Using cross-fitting makes it possible to obtain smaller standard errors for the DML estimator: The empirical distribution of the double machine learning estimator that is based on the cross-fitting Algorithm 2 (DML2) exhibits a more pronounced concentration around zero.

### *Inference on a structural parameter in key causal models*

We provide simulation results for double machine learning estimators in the presented key causal models in Figure 8. In a replication of the simulation example in Section 4.1, we show that the confidence intervals for the DML estimator in the partially linear regression model achieves an empirical coverage (= 0.952) close to the specified level of  $1 - \alpha = 0.95$ . The estimator is, again, based on a random forest learner. The corresponding results are presented in the top-left panel of Figure 8.

In a simulated example of a PLIV model, the DML confidence interval that is based on a lasso learner (`regr.cv_glmnet` of `mlr3`) achieves a coverage of 95.6%. The underlying data generating process is based on a setting considered in Chernozhukov, Hansen, and Spindler (2015a) with one instrumental variable. Moreover for simulations of the IRM model, we make use of a DGP of Belloni, Chernozhukov, Fernández-Val, and Hansen (2017). The DGP for the IIVM is inspired by a simulation run in Farbmacher, Guber, and Klaassen (2020). We present the formal DGPs in the appendix. To perform estimation of the nuisance functions in the

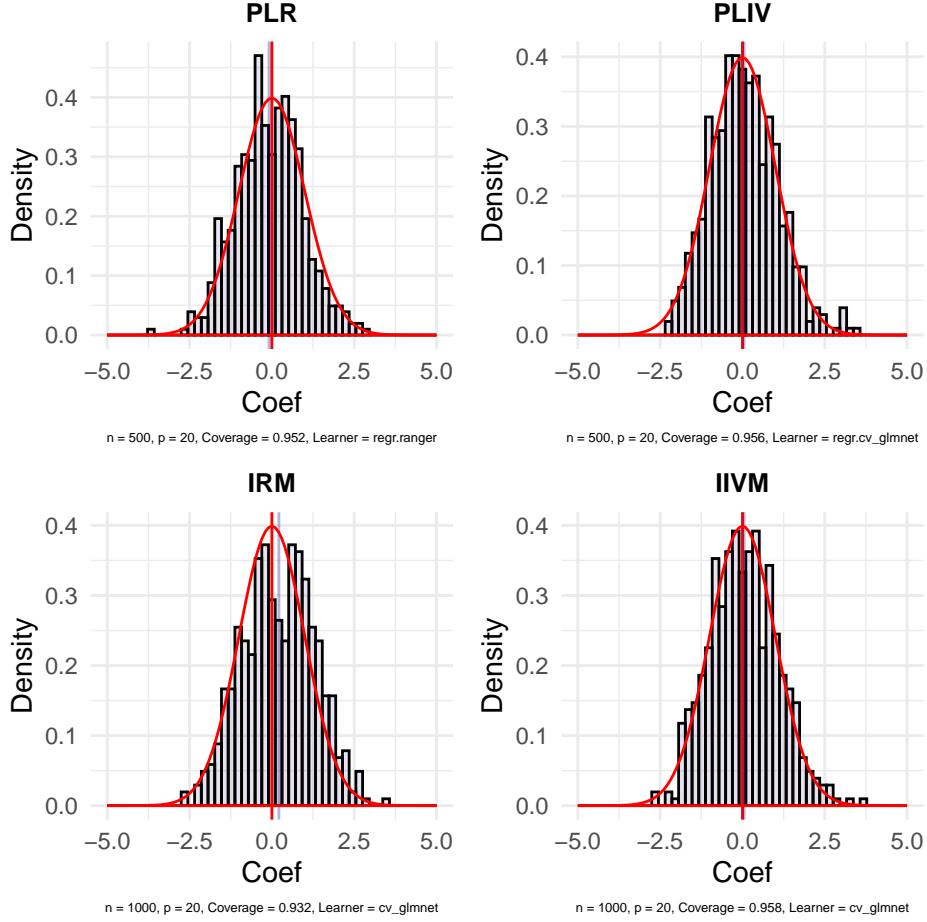


Figure 8: Histogram of double machine learning estimators in key causal models. The figure shows the histograms of the realizations of the DML estimators in the PLR (top left), PLIV (top right), IRM (bottom left), and IIVM (bottom right) as obtained in  $R = 500$  independent repetitions. Additional information on the data generating processes and implemented estimators are presented in the main text and the appendix.

interactive models, we employ the regression and classification predictors `regr.cv_glmnet` and `classif.cv_glmnet` as provided by the `mlr3` package. In all cases, we employ the cross-validated `lambda.min` choice of the penalty parameter with five folds, in other words, that  $\lambda$  value that minimizes the cross-validated mean squared error. Figure 8 shows that the empirical distribution of the centered estimators as obtained in finite sample settings is relatively well-approximated by a normal distribution. In all models the empirical coverage that is achieved by the constructed confidence bands is close to the nominal level.

### *Simultaneous inference*

To verify the finite-sample performance of the implemented methods for simultaneous inference, we perform a small simulation study in a regression setup which is similar as the one used in [Bach et al. \(2018\)](#). We would like to perform valid simultaneous inference on the



	CI	RW	Bonf.	Holm
FWER	0.08	0.11	0.08	0.10
Correct rejections	12.00	12.00	12.00	12.00

Table 1: Family-wise error rate (FWER) and average number of correct rejections in a simulation example. CI: Joint confidence interval, RW: Romano-Wolf stepdown correction, Bonf.: Bonferroni adjustment, Holm: Bonferroni-Holm correction.

coefficients  $\theta$  in the regression model

$$y_i = \beta_0 + d_i^\top \theta + \varepsilon_i, \quad i = 1, \dots, n,$$

with  $n = 1000$  and  $p_1 = 42$  regressors. The errors  $\varepsilon_i$  are normally distributed with  $\varepsilon_i \sim N(0, \sigma^2)$  and variance  $\sigma^2 = 3$ . The regressors  $d_i$  are generated by a joint normal distribution  $d_i \sim N(\mu, \Sigma)$  with  $\mu = \mathbf{0}$  and  $\Sigma_{j,k} = 0.5^{|j-k|}$ . The model is sparse in that only the first  $s = 12$  regressors have a non-zero effect on outcome  $y_i$ . The  $p_1$  coefficients  $\theta_1, \dots, \theta_{p_1}$  are generated as

$$\theta_j = \min \left\{ \frac{\theta^{\max}}{j^a}, \theta^{\min} \right\},$$

for  $j = 1, \dots, s$  with  $\theta^{\max} = 9$ ,  $\theta^{\min} = 0.75$ , and  $a = 0.99$ . All other coefficients have values exactly equal to 0. Estimation of the nuisance components has been performed by using the lasso as provided by `regr.cv_glmnet` in **mlr3**.

We report the empirical coverage as achieved by a joint  $(1 - \alpha)$ -confidence interval for all  $p_1 = 42$  coefficients and the realized family-wise error rate of the implemented  $p$  value adjustments in  $R = 500$  repetitions in Table 1. The finite sample performance of the Romano-Wolf stepdown procedure that is based on the multiplier bootstrap as well as the classical Bonferroni and Bonferroni-Holm correction are evaluated. Table 1 shows that all methods achieve an empirical FWER close to the specified level of  $\alpha = 0.1$ . In all cases, the double machine learning estimators reject all 12 false null hypotheses in every repetition.

## 8. Conclusion

In this paper, we provide an overview on the key ingredients and the major structure of the double/debiased machine learning framework as established in Chernozhukov *et al.* (2018) together with an overview on a collection of structural models. Moreover, we introduce the R package **DoubleML** that serves as an implementation of the double machine learning approach. A brief simulation study provides insights on the finite sample performance of the double machine learning estimator in the key causal models.

The structure of **DoubleML** is intended to be flexible with regard to the implemented structural models, the resampling scheme, the machine learning methods and the underlying algorithm, as well as the Neyman-orthogonal scores considered. By providing the R package **DoubleML** together with its Python twin (Bach *et al.* 2022), we hope to make double machine learning more accessible to users in practice. Finally, we would like to encourage users to add new structural models, scores and functionalities to the package.

## Acknowledgments

This work was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – Project Number 431701914.

## References

- Angrist JD, Pischke JS (2009). *Mostly Harmless Econometrics: An Empiricist’s Companion*. Princeton University Press.
- Athey S, Tibshirani J, Wager S (2019). “Generalized Random Forests.” *The Annals of Statistics*, **47**(2), 1148–1178. doi:10.1214/18-aos1709.
- Bach P, Chernozhukov V, Kurz MS, Spindler M (2022). “**DoubleML** – An Object-Oriented Implementation of Double Machine Learning in Python.” *Journal of Machine Learning Research*, **23**(53), 1–6. URL <http://jmlr.org/papers/v23/21-0862.html>.
- Bach P, Chernozhukov V, Kurz MS, Spindler M, Sven K (2024). *DoubleML: Double Machine Learning in R*. R package version 1.0.0, URL <https://CRAN.R-project.org/package=DoubleML>.
- Bach P, Chernozhukov V, Spindler M (2018). “Valid Simultaneous Inference in High-Dimensional Settings (with the **hdm** Package for R).” *arXiv 1809.04951*, arXiv.org E-Print Archive. doi:10.48550/arXiv.1809.04951.
- Battocchi K, Dillon E, Hei M, Lewis G, Oka P, Oprescu M, Syrgkanis V (2019). “**EconML**: A Python Package for ML-Based Heterogeneous Treatment Effects Estimation.” Python package version 0.14.1, URL <https://github.com/py-why/EconML>.
- Becker M, Binder M, Bischl B, Lang M, Pfisterer F, Reich NG, Richter J, Schratz P, Sonabend R (2021). “**mlr3** Book.” URL <https://mlr3book.mlr-org.com/>.
- Becker M, Lang M, Richter J, Bischl B, Schalk D (2023). *mlr3tuning: Tuning for mlr3*. R package version 0.19.0, URL <https://CRAN.R-project.org/package=mlr3tuning>.
- Belloni A, Chen D, Chernozhukov V, Hansen C (2012). “Sparse Models and Methods for Optimal Instruments with an Application to Eminent Domain.” *Econometrica*, **80**, 2369–2429. doi:10.3982/ecta9626.
- Belloni A, Chernozhukov V, Chetverikov D, Wei Y (2018). “Uniformly Valid Post-Regularization Confidence Regions for Many Functional Parameters in Z-Estimation Framework.” *The Annals of Statistics*, **46**(6B), 3643–3675. doi:10.1214/17-aos1671.
- Belloni A, Chernozhukov V, Fernández-Val I, Hansen C (2017). “Program Evaluation and Causal Inference with High-Dimensional Data.” *Econometrica*, **85**(1), 233–298. doi:10.3982/ecta12723.
- Belloni A, Chernozhukov V, Hansen C (2011). “Inference for High-Dimensional Sparse Econometric Models.” In D Acemoglu, M Arellano, E Dekel (eds.), *Advances in Economics and Econometrics – Tenth World Congress*, pp. 245–295. Cambridge University Press, New York. ArXiv version available at doi:10.48550/arXiv.1201.0220.

- Belloni A, Chernozhukov V, Kato K (2014a). “Uniform Post-Selection Inference for Least Absolute Deviation Regression and Other Z-Estimation Problems.” *Biometrika*, **102**(1), 77–94. doi:[10.1093/biomet/asu056](https://doi.org/10.1093/biomet/asu056).
- Belloni A, Chernozhukov V, Wang L (2014b). “Pivotal Estimation via Square-Root Lasso in Nonparametric Regression.” *The Annals of Statistics*, **42**(2), 757–788. doi:[10.1214/14-aos1204](https://doi.org/10.1214/14-aos1204).
- Belloni A, Chernozhukov V, Hansen C (2014c). “Inference on Treatment Effects after Selection among High-Dimensional Controls.” *The Review of Economic Studies*, **81**(2 (287)), 608–650. doi:[10.1093/restud/rdt044](https://doi.org/10.1093/restud/rdt044).
- Bia M, Huber M, Laffers L (2020). “Double Machine Learning for Sample Selection Models.” *arXiv 2012.00745*, arXiv.org E-Print Archive. doi:[10.48550/ARXIV.2012.00745](https://doi.org/10.48550/ARXIV.2012.00745).
- Bickel PJ, Klaassen CAJ, Ritov Y, Wellner JA (1993). *Efficient and Adaptive Estimation for Semiparametric Models*, volume 4. Johns Hopkins University Press Baltimore.
- Biliyas Y (2000). “Sequential Testing of Duration Data: The Case of the Pennsylvania ‘Reemployment Bonus’ Experiment.” *Journal of Applied Econometrics*, **15**(6), 575–594. doi:[10.1002/jae.579](https://doi.org/10.1002/jae.579).
- Bodory H, Huber M (2023). **causalweight**: *Estimation Methods for Causal Inference Based on Inverse Probability Weighting*. R package version 1.0.4, URL <https://CRAN.R-project.org/package=causalweight>.
- Breiman L (1996). “Stacked Regressions.” *Machine Learning*, **24**(1), 49–64. doi:[10.1007/bf00117832](https://doi.org/10.1007/bf00117832).
- Bühlmann P, Van de Geer S (2011). *Statistics for High-Dimensional Data: Methods, Theory and Applications*. Springer-Verlag. doi:[10.1007/978-3-642-20192-9](https://doi.org/10.1007/978-3-642-20192-9).
- Bühlmann P, Van de Geer S (2015). “High-Dimensional Inference in Misspecified Linear Models.” *Electronic Journal of Statistics*, **9**(1), 1449–1473. doi:[10.1214/15-ejs1041](https://doi.org/10.1214/15-ejs1041).
- Chang NC (2020). “Double/Debiased Machine Learning for Difference-in-Differences Models.” *The Econometrics Journal*, **23**(2), 177–191. doi:[10.1093/ectj/utaa001](https://doi.org/10.1093/ectj/utaa001).
- Chang W (2021). **R6**: *Encapsulated Classes with Reference Semantics*. R package version 2.5.1, URL <https://CRAN.R-project.org/package=R6>.
- 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. doi:[10.1111/ectj.12097](https://doi.org/10.1111/ectj.12097).
- Chernozhukov V, Chetverikov D, Kato K (2013). “Gaussian Approximations and Multiplier Bootstrap for Maxima of Sums of High-Dimensional Random Vectors.” *The Annals of Statistics*, **41**(6), 2786–2819. doi:[10.1214/13-aos1161](https://doi.org/10.1214/13-aos1161).
- Chernozhukov V, Chetverikov D, Kato K (2014). “Gaussian Approximation of Suprema of Empirical Processes.” *The Annals of Statistics*, **42**(4), 1564–1597. doi:[10.1214/14-aos1230](https://doi.org/10.1214/14-aos1230).

- Chernozhukov V, Hansen C, Spindler M (2015a). “Post-Selection and Post-Regularization Inference in Linear Models with Many Controls and Instruments.” *American Economic Review*, **105**(5), 486–90. doi:[10.1257/aer.p20151022](https://doi.org/10.1257/aer.p20151022).
- Chernozhukov V, Hansen C, Spindler M (2015b). “Valid Post-Selection and Post-Regularization Inference: An Elementary, General Approach.” *Annual Review of Economics*, **7**(1), 649–688. doi:[10.1146/annurev-economics-012315-015826](https://doi.org/10.1146/annurev-economics-012315-015826).
- Chernozhukov V, Hansen C, Spindler M (2016). “**hdm**: High-Dimensional Metrics.” *The R Journal*, **8**(2), 185–199. doi:[10.32614/RJ-2016-040](https://doi.org/10.32614/RJ-2016-040).
- Coyle J (2021). **tmle3**: *The Extensible TMLE Framework*. doi:[10.5281/zenodo.4603358](https://doi.org/10.5281/zenodo.4603358). R package version 0.2.0.
- Coyle J, Hejazi N, Malenica I, Phillips R, Sofrygin O (2021). **sl3**: *Pipelines for Machine Learning and Super Learning*. doi:[10.5281/zenodo.1342293](https://doi.org/10.5281/zenodo.1342293). R package version 1.4.4.
- Csárdi G, Hester J, Wickham H, Chang W, Morgan M, Tenenbaum D (2023). **remotes**: *R Package Installation from Remote Repositories, Including GitHub*. R package version 2.4.2.1, URL <https://CRAN.R-project.org/package=remotes>.
- Cunningham S (2021). *Causal Inference: The Mixtape*. Yale University Press.
- Dezeure R, Bühlmann P, Meier L, Meinshausen N (2015). “High-Dimensional Inference: Confidence Intervals,  $p$ -Values and R-Software **hdi**.” *Statistical Science*, **30**(4), 533–558. doi:[10.1214/15-sts527](https://doi.org/10.1214/15-sts527).
- Dowle M, Srinivasan A (2023). **data.table**: *Extension of data.frame*. R package version 1.14.8, URL <https://CRAN.R-project.org/package=data.table>.
- Dudoit S, Van der Laan MJ (2005). “Asymptotics of Cross-Validated Risk Estimation in Estimator Selection and Performance Assessment.” *Statistical Methodology*, **2**(2), 131–154. doi:[10.1016/j.stamet.2005.02.003](https://doi.org/10.1016/j.stamet.2005.02.003).
- Farbmacher H, Guber R, Klaassen S (2020). “Instrument Validity Tests with Causal Forests.” *Journal of Business & Economic Statistics*, pp. 1–10. doi:[10.1080/07350015.2020.1847122](https://doi.org/10.1080/07350015.2020.1847122).
- Farbmacher H, Huber M, Lafférs L, Langen H, Spindler M (2022). “Causal Mediation Analysis with Double Machine Learning.” *The Econometrics Journal*, **25**(2), 277–300. ISSN 1368-4221. doi:[10.1093/ectj/utac003](https://doi.org/10.1093/ectj/utac003). <https://academic.oup.com/ectj/article-pdf/25/2/277/43772863/utac003.pdf>.
- Farrell MH, Liang T, Misra S (2021). “Deep Neural Networks for Estimation and Inference.” *Econometrica*, **89**(1), 181–213. doi:[10.3982/ecta16901](https://doi.org/10.3982/ecta16901).
- Friedman J, Hastie T, Tibshirani R (2010). “Regularization Paths for Generalized Linear Models via Coordinate Descent.” *Journal of Statistical Software*, **33**(1), 1–22. doi:[10.18637/jss.v033.i01](https://doi.org/10.18637/jss.v033.i01).
- Gruber S, Van der Laan MJ (2012). “**tmle**: An R Package for Targeted Maximum Likelihood Estimation.” *Journal of Statistical Software*, **51**(13), 1–35. doi:[10.18637/jss.v051.i13](https://doi.org/10.18637/jss.v051.i13).

- Hasminskii RZ, Ibragimov IA (1978). “On the Nonparametric Estimation of Functionals.” In *Proceedings of the 2nd Prague Symposium on Asymptotic Statistics*, pp. 41–51.
- Imbens GW, Angrist JD (1994). “Identification and Estimation of Local Average Treatment Effects.” *Econometrica*, **62**(2), 467–475. doi:10.2307/2951620.
- Javanmard A, Montanari A (2014). “Hypothesis Testing in High-Dimensional Regression under the Gaussian Random Design Model: Asymptotic Theory.” *IEEE Transactions on Information Theory*, **60**(10), 6522–6554. doi:10.1109/tit.2014.2343629.
- Javanmard A, Montanari A (2018). “Debiasing the Lasso: Optimal Sample Size for Gaussian Designs.” *The Annals of Statistics*, **46**(6A), 2593–2622. doi:10.1214/17-aos1630.
- Klaassen CA (1987). “Consistent Estimation of the Influence Function of Locally Asymptotically Linear Estimators.” *The Annals of Statistics*, **15**(4), 1548–1562. doi:10.1214/aos/1176350609.
- Knaus MC (2021). “A Double Machine Learning Approach to Estimate the Effects of Musical Practice on Student’s Skills.” *Journal of the Royal Statistical Society A*, **184**(1), 282–300. doi:10.1111/rssa.12623.
- Knaus MC (2022). “Double Machine Learning-Based Programme Evaluation under Unconfoundedness.” *The Econometrics Journal*, **25**(3), 602–627. doi:10.1093/ectj/utac015.
- Kurz MS (2021). “Distributed Double Machine Learning with a Serverless Architecture.” In *Companion of the ACM/SPEC International Conference on Performance Engineering, ICPE ’21*, pp. 27–33. Association for Computing Machinery, New York. doi:10.1145/3447545.3451181.
- Lang M, Au Q, Coors S, Schratz P (2023a). **mlr3learners**: Recommended Learners for **mlr3**. R package version 0.5.6, URL <https://CRAN.R-project.org/package=mlr3learners>.
- Lang M, Binder M, Richter J, Schratz P, Pfisterer F, Coors S, Au Q, Casalicchio G, Kotthoff L, Bischl B (2019). “**mlr3**: A Modern Object-Oriented Machine Learning Framework in R.” *Journal of Open Source Software*, **4**(44), 1903. doi:10.21105/joss.01903.
- Lang M, Bischl B, Richter J, Sun X, Binder M (2023b). **paradox**: Define and Work with Parameter Spaces for Complex Algorithms. R package version 0.11.1, URL <https://CRAN.R-project.org/package=paradox>.
- Levit BY (1975). “On Efficiency of a Class of Non-Parametric Estimates.” *Teoriya Veroyatnostei i Ee Primeneniya*, **20**(4), 738–754. doi:10.1137/1120081.
- Newey W (1994). “The Asymptotic Variance of Semiparametric Estimators.” *Econometrica*, **62**(6), 1349–1382. doi:10.2307/2951752.
- Neyman J (1959). “Optimal Asymptotic Tests of Composite Hypotheses.” In U Grenander (ed.), *Probability and Statistics*, pp. 213–234. Almqvist & Wiksell.
- Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M, Duchesnay É (2011). “**Scikit-learn**: Machine Learning in Python.” *Journal*

- of Machine Learning Research*, **12**(85), 2825–2830. URL <http://jmlr.org/papers/v12/pedregosa11a.html>.
- Permutt T, Hebel JR (1989). “Simultaneous-Equation Estimation in a Clinical Trial of the Effect of Smoking on Birth Weight.” *Biometrics*, **45**, 619–622. doi:10.2307/2531503.
- Polley E, LeDell E, Kennedy C, Van der Laan MJ (2023). **SuperLearner**: *Super Learner Prediction*. R package version 2.0-28.1, URL <https://CRAN.R-project.org/package=SuperLearner>.
- Poterba JM, Venti SF, Wise DA (1994). “401(k) Plans and Tax-Deferred Saving.” *Studies in the Economics of Aging*, pp. 105–142. doi:10.3386/w4181.
- Poterba JM, Venti SF, Wise DA (1995). “Do 401(k) Contributions Crowd out Other Personal Saving?” *Journal of Public Economics*, **58**(1), 1–32. doi:10.1016/0047-2727(94)01462-w.
- R Core Team (2023). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- Robins JM, Rotnitzky A (1995). “Semiparametric Efficiency in Multivariate Regression Models with Missing Data.” *Journal of the American Statistical Association*, **90**(429), 122–129. doi:10.1080/01621459.1995.10476494.
- Robinson PM (1988). “Root-N-Consistent Semiparametric Regression.” *Econometrica*, **56**(4), 931–954. doi:10.2307/1912705.
- Romano JP, Wolf M (2005a). “Exact and Approximate Stepdown Methods for Multiple Hypothesis Testing.” *Journal of the American Statistical Association*, **100**(469), 94–108. doi:10.1198/016214504000000539.
- Romano JP, Wolf M (2005b). “Stepwise Multiple Testing as Formalized Data Snooping.” *Econometrica*, **73**(4), 1237–1282. doi:10.1111/j.1468-0262.2005.00615.x.
- Romano JP, Wolf M (2016). “Efficient Computation of Adjusted  $p$ -Values for Resampling-Based Stepdown Multiple Testing.” *Statistics & Probability Letters*, **113**, 38–40. doi:10.1016/j.spl.2016.02.012.
- Schick A (1986). “On Asymptotically Efficient Estimation in Semiparametric Models.” *The Annals of Statistics*, **14**(3), 1139–1151. doi:10.1214/aos/1176350055.
- Sonabend R, Schratz P, Fischer S (2023). **mlr3extralearners**: *Extra Learners for Use in mlr3*. R package version 0.7.1, URL <https://mlr3extralearners.ml-org.com/>.
- Tibshirani J, Athey S, Wager S (2023). **grf**: *Generalized Random Forests*. R package version 2.3.0, URL <https://CRAN.R-project.org/package=grf>.
- Van de Geer S, Bühlmann P, Ritov Y, Dezeure R (2014). “On Asymptotically Optimal Confidence Regions and Tests for High-Dimensional Models.” *The Annals of Statistics*, **42**(3), 1166–1202. doi:10.1214/14-aos1221.

- Van der Laan MJ, Coyle JR, Hejazi NS, Malenica I, Phillips RV, Hubbard AE (2022). “Targeted Learning in R: Causal Data Science with the **tlverse** Software Ecosystem.” URL <https://tlverse.org/tlverse-handbook/>.
- Van der Laan MJ, Polley EC, Hubbard AE (2007). “Super Learner.” *Statistical Applications in Genetics and Molecular Biology*, **6**(1). doi:10.2202/1544-6115.1309.
- Van der Laan MJ, Rose S (2011). *Targeted Learning: Causal Inference for Observational and Experimental Data*. Springer-Verlag. doi:10.1007/978-1-4419-9782-1.
- Van der Vaart AW (2000). *Asymptotic Statistics*, volume 3. Cambridge University Press.
- Van Rossum G, et al. (2011). *Python Programming Language*. URL <http://www.python.org/>.
- Wickham H (2019). *Advanced R*. CRC Press.
- Wolpert DH (1992). “Stacked Generalization.” *Neural Networks*, **5**(2), 241–259. doi:10.1016/s0893-6080(05)80023-1.
- Wright MN, Ziegler A (2017). “**ranger**: A Fast Implementation of Random Forests for High Dimensional Data in C++ and R.” *Journal of Statistical Software*, **77**(1), 1–17. doi:10.18637/jss.v077.i01.
- Wright PG (1928). *Tariff on Animal and Vegetable Oils*. Macmillan Company, New York.
- Zhang CH, Zhang SS (2014). “Confidence Intervals for Low Dimensional Parameters in High Dimensional Linear Models.” *Journal of the Royal Statistical Society B*, **76**(1), 217–242. doi:10.1111/rssb.12026.
- Zheng W, Van der Laan MJ (2011). “Cross-Validated Targeted Minimum-Loss-Based Estimation.” In *Targeted Learning*, pp. 459–474. Springer-Verlag.
- Zhong Y, Naimi A (2021). **AIPW**: *Augmented Inverse Probability Weighting (AIPW) for Binary Exposure*. R package version 0.6.3.1, URL <https://github.com/yqzhong7/AIPW>.

## A. Computation and infrastructure

The code in the paper has been executed with **DoubleML**, version 0.5.3.

The simulation study has been run on a x86\_64, darwin17.0 with macos Big Sur ... 10.16 system using R version 4.2.3 (2023-03-15). The following packages have been used for estimation: **DoubleML**, version 0.5.3, **data.table**, version 1.14.6, **mlr3**, version 0.14.1, **mlr3tuning**, version 0.17.2, **mlr3learners**, version 0.5.5, **glmnet**, version 4.1-6, **ranger**, version 0.14.1, **paradox**, version 0.11.0, **foreach**, version 1.5.2.

## B. Suppressed code output

### B.1. Pennsylvania Reemployment Data, Section 7

Load data as `data.table`.

```
R> library("DoubleML")
R> dt_bonus <- fetch_bonus(return_type = "data.table")
R> dt_bonus
```

	inuidur1	female	black	othrace	dep1	dep2	q2	q3	q4	q5	q6	age1t35	agegt54
1:	2.890372	0	0	0	0	1	0	0	0	1	0	0	0
2:	0.000000	0	0	0	0	0	0	0	0	1	0	0	0
3:	3.295837	0	0	0	0	0	0	0	1	0	0	0	0
4:	2.197225	0	0	0	0	0	0	1	0	0	0	1	0
5:	3.295837	0	0	0	1	0	0	0	0	1	0	0	1
---													
5095:	2.302585	0	0	0	0	0	0	1	0	0	0	1	0
5096:	1.386294	0	0	0	0	1	1	0	0	0	0	0	0
5097:	2.197225	0	0	0	0	1	1	0	0	0	0	1	0
5098:	1.386294	0	0	0	0	0	0	0	0	1	0	0	1
5099:	3.295837	0	0	0	0	0	0	0	1	0	0	0	1
		durable	lusrd	husd	tg								
1:		0	0	1	0								
2:		0	1	0	0								
3:		0	1	0	0								
4:		0	0	0	1								
5:		1	1	0	0								
---													
5095:		0	0	0	1								
5096:		0	0	0	1								
5097:		0	1	0	0								
5098:		0	0	0	1								
5099:		1	1	0	0								

```
R> obj_dml_data_bonus <- DoubleMLData$new(dt_bonus,
+   y_col = "inuidur1", d_cols = "tg",
+   x_cols = c("female", "black", "othrace", "dep1", "dep2", "q2", "q3",
+   "q4", "q5", "q6", "age1t35", "agegt54", "durable", "lusrd", "husd"))
```

Print data backend: Lists main attributes and methods of a 'DoubleMLData' object.



```
R> obj_dml_data_bonus
```

Print data set.

```
R> obj_dml_data_bonus$data
```

```

      inuidur1 female black othrace dep1 dep2 q2 q3 q4 q5 q6 agelt35 agegt54
1: 2.890372      0    0      0    0    1 0 0 0 1 0      0      0
2: 0.000000      0    0      0    0    0 0 0 0 1 0      0      0
3: 3.295837      0    0      0    0    0 0 0 1 0 0      0      0
4: 2.197225      0    0      0    0    0 0 1 0 0 0      1      0
5: 3.295837      0    0      0    1    0 0 0 0 1 0      0      1
---
5095: 2.302585      0    0      0    0    0 0 1 0 0 0      1      0
5096: 1.386294      0    0      0    0    1 1 0 0 0 0      0      0
5097: 2.197225      0    0      0    0    1 1 0 0 0 0      1      0
5098: 1.386294      0    0      0    0    0 0 0 0 1 0      0      1
5099: 3.295837      0    0      0    0    0 0 0 1 0 0      0      1
      durable lUSD husd tg
1:      0      0      1 0
2:      0      1      0 0
3:      0      1      0 0
4:      0      0      0 1
5:      1      1      0 0
---
5095:      0      0      0 1
5096:      0      0      0 1
5097:      0      1      0 0
5098:      0      0      0 1
5099:      1      1      0 0

```

```
R> learner_classif_m <- lrn("classif.ranger", num.trees = 500,
+   min.node.size = 2, max.depth = 5)
R> doubleml_irm_bonus <- DoubleMLIRM$new(obj_dml_data_bonus,
+   ml_g = learner_g, ml_m = learner_classif_m, score = "ATE",
+   dml_procedure = "dml1", n_folds = 5, n_rep = 1)
R> doubleml_irm_bonus
```

```
===== DoubleMLIRM Object =====
```

```
----- Data summary -----
```

```
Outcome variable: inuidur1
```

```
Treatment variable(s): tg
```

```
Covariates: female, black, othrace, dep1, dep2, q2, q3,
q4, q5, q6, agelt35, agegt54, durable, lUSD, husd
```

```
Instrument(s):
```

```
No. Observations: 5099
```

```
----- Score & algorithm -----
```

```
Score function: ATE
```

```
DML algorithm: dml1
```

```
----- Machine learner -----
```

```
ml_g: regr.ranger
ml_m: classif.ranger
```

```
----- Resampling -----
No. folds: 5
No. repeated sample splits: 1
Apply cross-fitting: TRUE
----- Fit summary -----
```

## B.2. Data backend with multiple treatment variables, Section 7.5

```
R> doubleml_data <- double_ml_data_from_data_frame(df, y_col = "y",
+   d_cols = c("X1", "X2", "X3", "X4", "X5", "X6", "X7", "X8", "X9", "X10"))
```

Set treatment variable d to X1.

```
R> doubleml_data
```

```
===== DoubleMLData Object =====
```

```
----- Data summary -----
Outcome variable: y
Treatment variable(s): X1, X2, X3, X4, X5, X6, X7, X8, X9, X10
Covariates: X11, X12, X13, X14, X15, X16, X17, X18, X19, X20, X21, X22, X23,
  X24, X25, X26, X27, X28, X29, X30, X31, X32, X33, X34, X35, X36, X37, X38,
  X39, X40, X41, X42, X43, X44, X45, X46, X47, X48, X49, X50, X51, X52, X53,
  X54, X55, X56, X57, X58, X59, X60, X61, X62, X63, X64, X65, X66, X67, X68,
  X69, X70, X71, X72, X73, X74, X75, X76, X77, X78, X79, X80, X81, X82, X83,
  X84, X85, X86, X87, X88, X89, X90, X91, X92, X93, X94, X95, X96, X97, X98,
  X99, X100
Instrument(s):
No. Observations: 500
```

## B.3. List of externally provided parameters, Section 7.6

```
R> str(doubleml_plr$params)
```

```
List of 2
 $ ml_l:List of 10
  ..$ X1 :List of 1
  .. ..$ lambda: num 0.09
  ..$ X2 :List of 1
  .. ..$ lambda: num 0.085
  ..$ X3 : NULL
  ..$ X4 : NULL
  ..$ X5 : NULL
  ..$ X6 : NULL
  ..$ X7 : NULL
```

```

..$ X8 : NULL
..$ X9 : NULL
..$ X10: NULL
$ ml_m:List of 10
..$ X1 :List of 1
.. ..$ lambda: num 0.1
..$ X2 :List of 1
.. ..$ lambda: num 0.095
..$ X3 : NULL
..$ X4 : NULL
..$ X5 : NULL
..$ X6 : NULL
..$ X7 : NULL
..$ X8 : NULL
..$ X9 : NULL
..$ X10: NULL

```

#### B.4. List of internally tuned parameters, Section 7.6

```
R> doubleml_plr$tuning_res$X1
```

```

$ml_1
$ml_1[[1]]
$ml_1[[1]]$tuning_result
$ml_1[[1]]$tuning_result[[1]]
$ml_1[[1]]$tuning_result[[1]]$tuning_result
  lambda learner_param_vals  x_domain regr.mse
1:    0.1      <list[2]> <list[1]> 10.53451

$ml_1[[1]]$tuning_result[[1]]$tuning_archive
  lambda regr.mse warnings errors runtime_learners
1:  0.100 10.53451      0      0              0.10
2:  0.095 10.60720      0      0              0.05
3:  0.085 10.76577      0      0              0.05
4:  0.055 11.32053      0      0              0.05
5:  0.060 11.21736      0      0              0.05
6:  0.050 11.42918      0      0              0.09
7:  0.075 10.93077      0      0              0.05
8:  0.065 11.11709      0      0              0.05
9:  0.080 10.84518      0      0              0.03
10: 0.070 11.02168      0      0              0.09
11: 0.090 10.68576      0      0              0.07

      uhash  x_domain      timestamp batch_nr
1: 2407e852-06a7-4756-ace6-42524bc37e34 <list[1]> 2023-01-31 14:49:37      1
2: 3a35f2c6-b78c-4416-89c9-e69158d2716b <list[1]> 2023-01-31 14:49:37      2
3: c78c69f3-3a70-4493-afec-121320689918 <list[1]> 2023-01-31 14:49:37      3
4: 3ffd8bcd-fd2a-46d0-b4b2-c9d0445fba2a <list[1]> 2023-01-31 14:49:37      4
5: ba275b12-edc5-4c79-8630-4c9c5285095a <list[1]> 2023-01-31 14:49:38      5
6: ff65786b-19fa-4393-9a9a-4627f07d2f9f <list[1]> 2023-01-31 14:49:38      6
7: 728bcbdef-cfad-4f65-875c-0428f1bc4339 <list[1]> 2023-01-31 14:49:38      7
8: 96dcc4eb-652a-4825-9481-2fa11b7b274a <list[1]> 2023-01-31 14:49:38      8

```

```

 9: e08f6536-0b71-4a82-919f-de35e1872c3f <list[1]> 2023-01-31 14:49:38      9
10: bb30cbc8-8324-441b-bd61-e75d9e22892c <list[1]> 2023-01-31 14:49:39     10
11: 6bac97d1-86cf-43ea-988a-4e7916a13460 <list[1]> 2023-01-31 14:49:39     11

```

```

$m1_1[[1]]$tuning_result[[1]]$params
NULL

```

```

$m1_1[[1]]$params
$m1_1[[1]]$params[[1]]
$m1_1[[1]]$params[[1]]$family
[1] "gaussian"

```

```

$m1_1[[1]]$params[[1]]$lambda
[1] 0.1

```

```

$m1_1$params
$m1_1$params[[1]]
$m1_1$params[[1]]$family
[1] "gaussian"

```

```

$m1_1$params[[1]]$lambda
[1] 0.1

```

```

$m1_m
$m1_m[[1]]
$m1_m[[1]]$tuning_result
$m1_m[[1]]$tuning_result[[1]]
$m1_m[[1]]$tuning_result[[1]]$tuning_result
  lambda learner_param_vals x_domain regr.mse
1:    0.1      <list[2]> <list[1]> 0.9794034

```

```

$m1_m[[1]]$tuning_result[[1]]$tuning_archive
  lambda regr.mse warnings errors runtime_learners
1: 0.090 0.9798230      0      0          0.04
2: 0.055 0.9971462      0      0          0.07
3: 0.075 0.9830963      0      0          0.05
4: 0.050 1.0045139      0      0          0.06
5: 0.100 0.9794034      0      0          0.06
6: 0.060 0.9907519      0      0          0.05
7: 0.065 0.9869171      0      0          0.06
8: 0.095 0.9797396      0      0          0.06
9: 0.085 0.9804282      0      0          0.04
10: 0.070 0.9848766      0      0          0.08
11: 0.080 0.9813190      0      0          0.06

      uhash x_domain      timestamp batch_nr
1: 06cd05b2-2daa-4600-a982-3d35037604a1 <list[1]> 2023-01-31 14:49:39      1
2: 4aa95e02-b7e3-49a0-b0ae-a5eea57cb686 <list[1]> 2023-01-31 14:49:39      2
3: 7a97b5c0-44cb-458a-9694-a4205cb8652e <list[1]> 2023-01-31 14:49:40      3
4: 190a5ba4-00ca-412b-a021-bfe92ba46375 <list[1]> 2023-01-31 14:49:40      4
5: e83d84d3-91a4-4969-8310-c6e289d3182c <list[1]> 2023-01-31 14:49:40      5

```

```

6: ab176601-1969-4424-8a14-df3c6cf74d84 <list[1]> 2023-01-31 14:49:40      6
7: 0c3d1269-412b-4734-8fe2-1b9bb7319479 <list[1]> 2023-01-31 14:49:40      7
8: a808a20c-b862-4665-a83e-85f7d6489536 <list[1]> 2023-01-31 14:49:41      8
9: f7cdb8f7-eb89-48fd-b295-5dae06f55cf7 <list[1]> 2023-01-31 14:49:41      9
10: 0a503598-db1a-48bd-8d75-6b3620a4a0e2 <list[1]> 2023-01-31 14:49:41     10
11: 47c429d6-7bee-4a77-875b-8d4cdf23f7ee <list[1]> 2023-01-31 14:49:41     11

```

```

$ml_m[[1]]$tuning_result[[1]]$params
NULL

```

```

$ml_m[[1]]$params
$ml_m[[1]]$params[[1]]
$ml_m[[1]]$params[[1]]$family
[1] "gaussian"

```

```

$ml_m[[1]]$params[[1]]$lambda
[1] 0.1

```

```

$ml_m$params
$ml_m$params[[1]]
$ml_m$params[[1]]$family
[1] "gaussian"

```

```

$ml_m$params[[1]]$lambda
[1] 0.1

```

The tuned parameters:

```
R> str(doubleml_plr$params)
```

```

List of 2
 $ ml_1:List of 10
  ..$ X1 :List of 2
  .. ..$ family: chr "gaussian"
  .. ..$ lambda: num 0.1
  ..$ X2 :List of 2
  .. ..$ family: chr "gaussian"
  .. ..$ lambda: num 0.1
  ..$ X3 :List of 2
  .. ..$ family: chr "gaussian"
  .. ..$ lambda: num 0.1
  ..$ X4 :List of 2
  .. ..$ family: chr "gaussian"
  .. ..$ lambda: num 0.09
  ..$ X5 :List of 2
  .. ..$ family: chr "gaussian"
  .. ..$ lambda: num 0.07
  ..$ X6 :List of 2
  .. ..$ family: chr "gaussian"
  .. ..$ lambda: num 0.085
  ..$ X7 :List of 2

```

```
.. ..$ family: chr "gaussian"
.. ..$ lambda: num 0.085
..$ X8 :List of 2
.. ..$ family: chr "gaussian"
.. ..$ lambda: num 0.08
..$ X9 :List of 2
.. ..$ family: chr "gaussian"
.. ..$ lambda: num 0.09
..$ X10:List of 2
.. ..$ family: chr "gaussian"
.. ..$ lambda: num 0.075
$ ml_m:List of 10
..$ X1 :List of 2
.. ..$ family: chr "gaussian"
.. ..$ lambda: num 0.1
..$ X2 :List of 2
.. ..$ family: chr "gaussian"
.. ..$ lambda: num 0.095
..$ X3 :List of 2
.. ..$ family: chr "gaussian"
.. ..$ lambda: num 0.095
..$ X4 :List of 2
.. ..$ family: chr "gaussian"
.. ..$ lambda: num 0.095
..$ X5 :List of 2
.. ..$ family: chr "gaussian"
.. ..$ lambda: num 0.1
..$ X6 :List of 2
.. ..$ family: chr "gaussian"
.. ..$ lambda: num 0.1
..$ X7 :List of 2
.. ..$ family: chr "gaussian"
.. ..$ lambda: num 0.1
..$ X8 :List of 2
.. ..$ family: chr "gaussian"
.. ..$ lambda: num 0.1
..$ X9 :List of 2
.. ..$ family: chr "gaussian"
.. ..$ lambda: num 0.1
..$ X10:List of 2
.. ..$ family: chr "gaussian"
.. ..$ lambda: num 0.1
```

## C. Data generating processes, simulation study

### C.1. Data generating process for PLIV simulation

The DGP is based on Chernozhukov *et al.* (2015a) and defined as

$$\begin{aligned} z_i &= \Pi x_i + \zeta_i, \\ d_i &= x_i^\top \gamma + z_i^\top \delta + u_i, \\ y_i &= \alpha d_i + x_i^\top \beta + \varepsilon_i, \end{aligned}$$

with

$$\begin{pmatrix} \varepsilon_i \\ u_i \\ \zeta_i \\ x_i \end{pmatrix} \sim \mathcal{N} \left( 0, \begin{pmatrix} 1 & 0.6 & 0 & 0 \\ 0.6 & 1 & 0 & 0 \\ 0 & 0 & 0.25 I_{p_n^z} & 0 \\ 0 & 0 & 0 & \Sigma \end{pmatrix} \right)$$

where  $\Sigma$  is a  $p_n^x \times p_n^x$  matrix with entries  $\Sigma_{kj} = 0.5^{|k-j|}$  and  $I_{p_n^z}$  is an identity matrix with dimension  $p_n^z \times p_n^z$ .  $\beta = \gamma$  is a  $p_n^x$ -vector with entries  $\beta = \frac{1}{j^2}$  and  $\Pi = (I_{p_n^z}, 0_{p_n^z \times (p_n^x - p_n^z)})$ . In the simulation example, we have one instrument, i.e.,  $p_n^z = 1$  and  $p_n^x = 20$  regressors  $x_i$ . In the simulation study, data sets with  $n = 500$  observations are generated in  $R = 500$  independent repetitions.

### C.2. Data generating process for IRM simulation

The DGP is based on a simulation study in Belloni *et al.* (2017) and defined as

$$\begin{aligned} d_i &= 1 \left\{ \frac{\exp(c_d x_i^\top \beta)}{1 + \exp(c_d x_i^\top \beta)} > v_i \right\}, \quad v_i \sim \mathcal{U}(0, 1), \\ y_i &= \theta d_i + c_y x_i^\top \beta d_i + \zeta_i, \quad \zeta_i \sim \mathcal{N}(0, 1), \end{aligned}$$

with covariates  $x_i \sim \mathcal{N}(0, \Sigma)$  where  $\Sigma$  is a matrix with entries  $\Sigma_{kj} = 0.5^{|k-j|}$ .  $\beta$  is a  $p_x$ -dimensional vector with entries  $\beta_j = \frac{1}{j^2}$  and the constants  $c_y$  and  $c_d$  are determined as

$$c_y = \sqrt{\frac{R_y^2}{(1 - R_y^2) \beta^\top \Sigma \beta}}, \quad c_d = \sqrt{\frac{(\pi^2/3) R_d^2}{(1 - R_d^2) \beta^\top \Sigma \beta}}.$$

We set the values of  $R_y^2 = 0.5$  and  $R_d^2 = 0.5$  and consider a setting with  $n = 1000$  and  $p = 20$ . Data generation and estimation have been performed in  $R = 500$  independent replications.

### C.3. Data generating process for IIVM simulation

The DGP is defined as

$$\begin{aligned} d_i &= 1 \{ \alpha_x Z + v_i > 0 \}, \\ y_i &= \theta d_i + x_i^\top \beta + u_i, \end{aligned}$$

with  $Z \sim \text{Bernoulli}(0.5)$  and

$$\begin{pmatrix} u_i \\ v_i \end{pmatrix} \sim \mathcal{N}\left(0, \begin{pmatrix} 1 & 0.3 \\ 0.3 & 1 \end{pmatrix}\right).$$

The covariates are drawn from a multivariate normal distribution with  $x_i \sim \mathcal{N}(0, \Sigma)$  with entries of the matrix  $\Sigma$  being  $\Sigma_{kj} = 0.5^{|j-k|}$  and  $\beta$  being a  $p_x$ -dimensional vector with  $\beta_j = \frac{1}{\beta^2}$ . The data generating process is inspired by a process used in a simulation in [Farbmacher et al. \(2020\)](#). In the simulation study, data sets with  $n = 1000$  observations and  $p_x = 20$  confounding variables  $x_i$  have been generated in  $R = 500$  independent repetitions.

### Affiliation:

Philipp Bach, Martin Spindler, Sven Klaassen  
 University of Hamburg  
 Chair of Statistics  
 Moorweidenstr. 18  
 20148 Hamburg, Germany  
 E-mail: [philipp.bach@uni-hamburg.de](mailto:philipp.bach@uni-hamburg.de), [martin.spindler@uni-hamburg.de](mailto:martin.spindler@uni-hamburg.de),  
[sven.klaassen@uni-hamburg.de](mailto:sven.klaassen@uni-hamburg.de)

Malte S. Kurz  
 Technical University of Munich  
 TUM School of Management  
 Arcisstr. 21  
 80333 Munich, Germany  
 E-mail: [malte.kurz@tum.de](mailto:malte.kurz@tum.de)

Victor Chernozhukov  
 Massachusetts Institute of Technology  
 Department of Economics and Center for Statistics and Data Science  
 50 Memorial Drive  
 Cambridge, MA 02139, United States of America  
 E-mail: [vchern@mit.edu](mailto:vchern@mit.edu)