

IZA DP No. 7972

**Treatment Evaluation with Multiple Outcome Periods
under Endogeneity and Attrition**

Markus Frölich
Martin Huber

February 2014

Treatment Evaluation with Multiple Outcome Periods under Endogeneity and Attrition

Markus Frölich

*Universität Mannheim
and IZA*

Martin Huber

Universität St.Gallen

Discussion Paper No. 7972
February 2014

IZA

P.O. Box 7240
53072 Bonn
Germany

Phone: +49-228-3894-0

Fax: +49-228-3894-180

E-mail: iza@iza.org

Any opinions expressed here are those of the author(s) and not those of IZA. Research published in this series may include views on policy, but the institute itself takes no institutional policy positions. The IZA research network is committed to the IZA Guiding Principles of Research Integrity.

The Institute for the Study of Labor (IZA) in Bonn is a local and virtual international research center and a place of communication between science, politics and business. IZA is an independent nonprofit organization supported by Deutsche Post Foundation. The center is associated with the University of Bonn and offers a stimulating research environment through its international network, workshops and conferences, data service, project support, research visits and doctoral program. IZA engages in (i) original and internationally competitive research in all fields of labor economics, (ii) development of policy concepts, and (iii) dissemination of research results and concepts to the interested public.

IZA Discussion Papers often represent preliminary work and are circulated to encourage discussion. Citation of such a paper should account for its provisional character. A revised version may be available directly from the author.

ABSTRACT

Treatment Evaluation with Multiple Outcome Periods under Endogeneity and Attrition^{*}

This paper develops a nonparametric methodology for treatment evaluation with multiple outcome periods under treatment endogeneity and missing outcomes. We use instrumental variables, pre-treatment characteristics, and short-term (or intermediate) outcomes to identify the average treatment effect on the outcomes of compliers (the subpopulation whose treatment reacts on the instrument) in multiple periods based on inverse probability weighting. Treatment selection and attrition may depend on both observed characteristics and the unobservable compliance type, which is possibly related to unobserved factors. We also provide a simulation study and apply our methods to the evaluation of a policy intervention targeting college achievement, where we find that controlling for attrition considerably affects the effect estimates.

JEL Classification: C14, C21, C23, C24, C26

Keywords: treatment effect, attrition, endogeneity, panel data, weighting

Corresponding author:

Markus Frölich
Chair of Econometrics
Universität Mannheim
L7, 3-5
68131 Mannheim
Germany
E-mail: froelich@uni-mannheim.de

^{*} We have benefited from comments by Michael Lechner, Blaise Melly, Conny Wunsch, the editor and associate editor, and three anonymous referees. Frölich acknowledges financial support from the Research Center (SFB) 884 "Political Economy of Reforms" Project B5, funded by the German Research Foundation (DFG). Huber acknowledges financial support from the Swiss National Science Foundation grant PBSGP1_138770.

1 Introduction

We develop a nonparametric methodology for evaluating the effect of an endogenous binary variable (referred to as treatment) in multiple outcomes periods where some outcomes are missing non-randomly due to non-response and attrition (e.g. survey non-response or truncation by death). Our identification strategy exploits an instrument (to control for treatment endogeneity), baseline covariates, and short-term (or intermediate) post-treatment variables to tackle the dynamic nature of the attrition problem. This in principle allows us to estimate the treatment effects also in later periods where the attrition problem is typically particularly severe.

The proposed methods appear important in the light of two fundamental trends that are currently observed in applied research in social sciences: First, the increasing use of randomized experiments and second, a growing interest in medium to long-term treatment effects of interventions, in order to see whether effects are sustainable. Even randomized experiments, which are frequently regarded as the gold standard for causal inference, are often plagued by imperfections such as noncompliance with treatment assignment and outcome attrition due to loss to follow-up. The noncompliance issue can be solved if it can be plausibly assumed that random treatment assignment provides a credible instrument for (endogenous) treatment take-up. While this is common practice for the identification of complier average causal effects (CACE) (also known as local average treatment effects, LATE) in experiments, see Imbens and Angrist (1994) and Angrist, Imbens, and Rubin (1996), our approach also tackles the attrition problem. The latter appears particularly relevant when noting the increasing importance of long-term evaluations of policy interventions, as e.g. in the assessment of active labor market policies, e.g. Lechner, Miquel, and Wunsch (2011), or of educational interventions, e.g. Angrist, Bettinger, and Kremer (2006).

To see the contribution of this paper, it appears useful to review previously suggested approaches to correct for attrition. The very common missing at random (MAR) restriction assumes non-response or attrition to be conditionally ignorable (i.e., independent of the potential outcomes) given observed characteristics, see for instance Rubin (1976), Little and Rubin (1987), Robins, Rotnitzky, and Zhao (1994), Robins, Rotnitzky, and Zhao (1995), Carroll, Ruppert, and Stefanski (1995), Shah, Laird, and Schoenfeld (1997), Fitzgerald, Gottschalk, and Moffitt (1998), and Abowd, Crepon, and Kramarz (2001). Frangakis and Rubin (1999) suggest a relaxation of MAR in experiments which they call latent ignorability (LI). Non-response is assumed to be ignorable conditional on observed characteristics *and* the latent (compliance) type, characterizing

how an individual's treatment state reacts on some instrument. See Barnard, Frangakis, Hill, and Rubin (2003), Frangakis, Brookmeyer, Varadhan, Safaeian, Vlahov, and Strathdee (2004), and Mealli, Imbens, Ferro, and Biggeri (2004) for related applications.

Approaches other than MAR and LI, permitting attrition to be related to unobservables in a general way, are referred to as non-ignorable non-response models. The earlier work, e.g. Heckman (1976), Hausman and Wise (1979), Bollinger and David (2001), and Chen, Wong, Dominik, and Steiner (2000), focussed on fully parameterized maximum likelihood estimation with identification often achieved only via functional form restrictions, see Little (1995) for an intuitive example. Instrumental variables for non-response and attrition offer an additional source of identification, see DiNardo, McCrary, and Sanbonmatsu (2006) for an application in an experimental context. In particular, such models allow for non-parametric identification and more flexible estimation, including the series regression approach of Das, Newey, and Vella (2003) and inverse probability weighting based on instruments for attrition as outlined in Huber (2012, 2013). While the standard framework consists of just one follow-up period, panel data sample selection models as suggested by Kyriazidou (1997, 2001) can be used to consider multiple periods as in this paper. In addition to dynamic attrition, Semykina and Wooldridge (2006) even allow for endogenous regressors, given that sufficiently many instruments to control for attrition and endogeneity are available.

An alternative to the assumptions discussed so far are methods that do not require a fully specified model for attrition, however, at the cost of sacrificing point identification. E.g., building on the partial identification literature (Robins, 1989, Manski, 1989, 1990), Zhang and Rubin (2003), Zhang, Rubin, and Mealli (2008), Imai (2008), and Lee (2009), among others, bound treatment effects in the presence of non-response under comparably mild restrictions. Another approach is multiple imputation of missing values, which goes back to Rubin (1977, 1978). Based on Bayesian techniques, multiple attrition models are used to impute multiple sets of plausible values for the missing data in order to obtain a probability interval for the parameter of interest. Finally, Rotnitzky, Robins, and Scharfstein (1998), Scharfstein, Rotnitzky, and Robins (1999), and Xie and Qian (2012) (who even allow for non-monotone non-response), among others, propose sensitivity checks for violations of MAR related to unobservables by varying the nuisance term causing non-ignorable attrition over a relevant range to examine the robustness of the results. By not considering arguably implausible attrition mechanisms, this approach likely yields more informative results than Manski-style bounds analysis and therefore provides a middle ground

between the latter and point identification.

In this paper, we propose a new nonparametric approach for point identification of the average causal effect on the compliers (those who are responsive to the instrument). We rely on pre-treatment covariates and (endogenous) post-treatment variables to control for attrition in a panel data framework as well as a single instrument (e.g., random assignment in an experiment) to tackle treatment endogeneity. (We only require a *single* instrument, which is important because instrumental variables are often hard to find in applications.) Our method for the evaluation of binary treatments provides three improvements compared to standard MAR. Firstly, we do not control for pre-treatment covariates only. That would ignore information about the intermediate variables, which presumably are important predictors of non-response in many empirical contexts. Secondly, we allow for treatment endogeneity which has rarely been considered under MAR. Exceptions are Yau and Little (2001) and Ding and Lehrer (2010), who, however, rely on considerably stronger functional form assumptions and in the latter case, on a difference-in-difference strategy rather than an instrument. Thirdly, in our main identification theorem, we develop a panel data extension of LI by permitting that attrition does not only depend on observables but also on the latent types.

It is also interesting to compare our framework to the literature on dynamic treatment regimes, e.g. Robins, Greenland, and Hu (1999), Murphy, van der Laan, and Robins (2001), and Lok, Gill, van der Vaart, and Robins (2004). If one were to consider attrition as a dynamic treatment regime, those methods could be adjusted to our situation. However, they are all based on a type of dynamic ignorability condition, which would correspond to a MAR assumption in our context. In contrast, we also allow for selection on the latent types and make use of an instrumental variable to overcome the endogeneity problems.

Our framework is also more general than the original LI assumption of Frangakis and Rubin (1999). Firstly, we permit two-sided noncompliance (i.e. the existence of never takers, who are never treated irrespective of the instrument, and of always takers, who are always treated) and extend LI to conditional LI given observables. Secondly, we consider multiple periods under comparably weak assumptions, whereas the literature conventionally imposes more structure and assesses only one outcome period, see for instance Peng, Little, and Raghunathan (2004). Note, however, that the identification problem considered in this paper is distinct from non-ignorable non-response and panel data sample selection models. I.e., we assume that conditional on

observed characteristics *and* the latent type, there are no further unobservables that are jointly related to attrition and the potential outcomes. Therefore, we do *not* require any additional instruments for non-response, which are typically hard to find in applications, see the discussion in Fitzgerald, Gottschalk, and Moffitt (1998). All in all, the methods proposed in this paper use less severe functional form and/or identifying assumptions than many non-response models invoked in recent empirical applications, see the examples in Preisser, Galecki, Lohman, and Wagenknecht (2000), Mattei and Mealli (2007), Shepherd, Redman, and Ankerst (2008), Zhang, Rubin, and Mealli (2009), Frumento, Mealli, Pacini, and Rubin (2012), and Wang, Rotnitzky, Lin, Millikan, and Thall (2012).

The remainder of this paper is organized as follows. Section 2 introduces a treatment effect model with endogeneity and multiple outcome periods and shows nonparametric identification under two distinct forms of attrition. For the ease of exposition, only two outcome periods are considered in the main text. A simulation study is provided in Section 3. Section 4 presents an application to a policy intervention aiming to increase college achievement previously analyzed by Angrist, Lang, and Oreopoulos (2009). Section 5 concludes. The (separate) online appendix presents identification in the more general case with several outcome periods along with the identification proofs, discusses the implications of our identifying assumptions in a parametric benchmark model, provides nonparametric and \sqrt{n} -consistent estimators based on kernel regression along with the proofs of their asymptotic properties, and includes an extended range of simulation studies.

2 Model and identification

Suppose we are interested in estimating the treatment effect of a binary variable $D \in \{0, 1\}$ on an outcome Y_t , where the subscript t denotes the period ($t = 1, 2, 3, \dots$) after the start of the treatment. All variables observed *prior to* the treatment are indexed by period zero and are denoted as X_0 . The *potential* outcomes Y_t^1 and Y_t^0 are the outcomes that would have been realized if D had been set to 1 or 0, respectively, by external intervention. (To avoid confusion between subscripts and superscripts we sometimes write $Y_{t=1}^0$ instead of Y_1^0 when referring to a specific time period.) In our nonparametric identification framework, two major issues have to be dealt with: endogenous treatment selection and missing outcome data due to attrition or non-response. The indicator variable R_t will denote whether in time period t outcome data is observed ($R_t = 1$) or missing

($R_t = 0$). We assume that information on the treatment D and baseline covariates X_0 is available for all individuals, but that individuals may not respond or drop out at follow-up data collection. In most applications, non-response increases at later follow-up periods.

2.1 Treatment endogeneity without attrition

Consider first the case without missing data. Imbens and Angrist (1994) and Angrist, Imbens, and Rubin (1996) have shown that in the presence of an instrumental variable (denoted by Z) satisfying particular assumptions, treatment effects are nonparametrically identified for a subset of the population, the so-called compliers. Adhering to their terminology, let D_i^z denote the *potential* treatment status of some individual i if Z_i were hypothetically set to z . For ease of exposition we will focus on a binary Z , which often occurs in experiments, even though the framework could be extended easily to non-binary discrete instruments, see e.g. Frölich (2007). The two binary potential treatment states D_i^0 and D_i^1 partition the population into four different types of individuals according to treatment behavior: the always takers (a) who are treated irrespective of the instrument ($D_i^1 = 1, D_i^0 = 1$), the never takers (n) who are never treated ($D_i^1 = 0, D_i^0 = 0$), the compliers (c) who only attend treatment if the instrument takes the value one ($D_i^1 = 1, D_i^0 = 0$), and the defiers (d) who only attend treatment if the instrument takes the value zero ($D_i^1 = 0, D_i^0 = 1$). As shortcut notation we will henceforth use \mathcal{T}_i for ‘type’ with $\mathcal{T}_i \in \{a, n, c, d\}$. Note that the type of any individual is only partially observed, i.e. *latent*, because the observed D and Z do not uniquely determine \mathcal{T} , as discussed in the appendix.

Abadie (2003) shows the nonparametric identification of the CACE (or LATE)

$$E[Y_t^1 - Y_t^0 | \mathcal{T} = c],$$

i.e. the effect for the compliers, under conditions implying conditional validity of the instrument given observed baseline characteristics, which we denote by X_0 :

$$\begin{aligned} \{Y_t^d, \mathcal{T}\} \perp\!\!\!\perp Z | X_0 & \quad \text{for } d \in \{0, 1\}, \\ \Pr(\mathcal{T} = d | X_0) = 0 & \quad \Pr(\mathcal{T} = c | X_0) > 0. \end{aligned}$$

The first line assumes independence between the instrument and the type and potential outcomes, conditional on X_0 . (Note that ‘ $\perp\!\!\!\perp$ ’ denotes statistical independence). It thus assumes random assignment of Z and an exclusion restriction with respect to the potential outcomes for given values

of the baseline covariates X_0 . The second line states that the treatment is (weakly) monotonic in the instrument conditional on X_0 so that defiers are ruled out and compliers do exist.

In the subsequent sections, we extend the CACE framework to allow for missing values in the outcome variables Y_t . We focus on the case of attrition (i.e. missingness as an absorbing state), which is the most frequent concern in empirical applications, particularly in impact evaluation. However, our approach does also permit intermittent missingness, implying that intermediate outcomes are missing while later ones are observed, but in this case does not exploit the information from later waves. With this respect, it is interesting to note that several contributions considering parametric missing data models distinguish explicitly between attrition and intermittent missingness, see for instance Xie and Qian (2012). In those approaches, however, one either has to additionally model the re-entry process after non-response or specify attrition and intermittent missingness as two separate processes. Under additional assumptions, also our nonparametric approach could use information from the re-entrants in order to permit more precise estimates (given that re-entry occurs sufficiently often), but the identification expressions and estimators would become less tractable. Since we aim at imposing as few restrictions as possible and do not make use of additional instruments (other than for treatment), which are often not available in applications, we therefore only model the non-response process and ignore any information after the first non-response. Hence, we do permit that individuals have missing data in only one or several waves and then re-enter the panel after periods of non-observability, but we do not exploit this information.

In the following, we denote by X_t the observed characteristics for any $t > 0$, i.e. *after* treatment. Note that X_t usually also contains the outcome Y_t . In contrast to X_0 , these variables X_t may possibly already be causally affected by the treatment, and we refer to them as (endogenous) post-treatment characteristics. (Note that whereas X_0 is permitted to be endogenous in the sense of Frölich (2008), i.e. that X_0 may be correlated with baseline unobservables, X_0 is not permitted to be causally influenced by treatment D , e.g. due to anticipation.) Furthermore, define $\underline{X}_t = \{X_1, \dots, X_t\}$ to be the history of the characteristics up to time t , where we do *not* include X_0 here in order to make the distinction between pre-treatment and endogenous post-treatment variables explicit. Accordingly, X_t^d and \underline{X}_t^d denote the *potential* values of the characteristics and of their history, respectively, at time t , if the treatment had been set to d by external intervention. Furthermore, let R_t be the response indicator in period t . I.e., X_t and Y_t are only observed if $R_t = 1$. Our setup permits that R_1 is zero for some individuals, such that outcome data is completely miss-

ing for those subjects. The history of response indicators over the post-treatment periods up to t is denoted by $\underline{R}_t = \{R_1, \dots, R_t\}$. The *potential* values of response and the response history are denoted by R_t^d and \underline{R}_t^d , respectively.

The occurrence of attrition and non-response may have many reasons. In the simplest and least realistic case, it is only triggered by random events happening after treatment, such that outcomes are missing completely at random (MCAR), see e.g. Rubin (1976) and Heitjan and Basu (1996). However, it is more likely that attrition depends also on observed and/or unobserved characteristics of the individuals. In particular, attrition may depend on Y_{t-1} , which is an endogenous variable that has been causally affected by the treatment. In addition, attrition could also be directly causally affected by the treatment itself, e.g. due to side effects or adverse events of a drug treatment in a medical intervention. Finally, attrition could also be caused directly by the instrumental variable Z .

Our identification strategy requires us to restrict the missing data process in two ways: First, we assume that non-response in time t is not *simultaneously* related to the outcome variable in time t . This implies that while any variables measured in the past may trigger non-response today, current and future values of the outcome variable are *not* permitted to do so. Non-response is thus considered to be predetermined. Second, we need to impose some restrictions on the relationship between the instrument and non-response. In the following two subsections, we will discuss two different identification assumptions. The first approach permits non-response to depend on unobservables, but requires it to be ignorable given the observed characteristics *and* the latent types (*conditional LI*). The second approach assumes that non-response is missing at random (MAR) given the observed pre- and post-treatment characteristics. While the first setup appears to be more general in most applications than the second one, they are not strictly nested. I.e. while the first approach is less restrictive with respect to the non-response process, the second one imposes weaker (albeit only mildly weaker) assumptions on the instrument.

Our analysis covers four cases. First, it includes randomized experiments with full compliance. Then, the exclusion restriction is valid with X_0 being the empty set (i.e. not controlling for any covariates) and using D as its own instrument, i.e. defining $Z_i \equiv D_i$. Second, under random assignment but imperfect compliance, we may use the randomization Z as an instrumental variable for the actual treatment receipt D . If the randomization probability is the same for everyone, X_0 may again be the empty set. Third, the framework also includes observational studies, where the

instrumental variable assumption is often plausible only after conditioning on some variables X_0 , see Abadie (2003), Tan (2006), and Frölich (2007). Finally, when controlling for some X_0 and using D as its own instrument, i.e. defining $Z_i \equiv D_i$, we impose what is referred to in the literature as the selection on observables, unconfoundedness, ignorable assignment, or conditional independence assumption, see for instance Rosenbaum and Rubin (1983), Lechner (1999), and Imbens (2004). Hence, although we focus on instrument-based identification, our identification results are also directly applicable to the selection on observables framework with missing outcome data.

2.2 Non-response under conditional latent ignorability

This section presents the identifying assumptions for the case of *conditional latent ignorability*. We permit that the response process at time R_t is related to all observed variables in the past *and* that it is a function of the latent type \mathcal{T} . Hence, the response process is supposed to be predetermined, which means that past values of the outcomes and further observed characteristics may affect the response behavior today. However, conditional on these past values, the instrument and conditional on the latent type, current and future outcomes must be independent of non-response in period t . This is, for instance, different to Xie and Qian (2012), who permit response and contemporaneous outcomes to be related and propose various sensitivity checks. Assumption 1 formalizes predetermined non-response under conditional LI. As already mentioned, \mathbf{X}_{t-1} may contain both intermediate outcomes \mathbf{Y}_{t-1} as well as other observed characteristics.

Assumption 1: Predetermined non-response

$$Y_{t+s} \perp\!\!\!\perp R_t | X_0, \mathbf{X}_{t-1}, \mathbf{R}_{t-1}, Z, \mathcal{T}, \quad \text{for } s \geq 0. \quad (1)$$

The plausibility of predetermined non-response (not related to contemporaneous outcomes) needs to be judged in the light of the application at hand. Some statistical support that this may be an empirically relevant case comes from Hirano, Imbens, Ridder, and Rubin (2001), who provide conditions implied by non-response related to (i) past information and (ii) contemporaneous outcomes that can be tested if a refreshment sample is available. Applying their test to a Dutch household survey, they reject attrition related to contemporaneous outcomes, but do not reject predetermined non-response at any conventional level. Our assumption may for instance appear plausible in the context of educational outcomes, where Y_t denotes a measure of cognitive skills (e.g. test scores or grades) at the end of some academic year t and R_t is

an indicator for (not) having dropped out of school. Predetermined non-response is (closely) satisfied if individuals decide to remain in or leave education (mainly) based on their academic performance in the previous academic year, Y_{t-1} , so that the drop-out decision R_t is taken shortly after that, e.g. during or at the end of summer vacation.

In addition to Assumption 1, we invoke exclusion, monotonicity, and common support restrictions, as stated in Assumptions 2 and 3. The latter are similar to Abadie (2003), apart from that we have to strengthen the instrumental exclusion restriction for the always and never takers.

Assumption 2: Exclusion restriction: For $d \in \{0, 1\}$

$$\begin{aligned} (Y_t, \mathbf{X}_{t-1}, \mathbf{R}_{t-1}) &\perp\!\!\!\perp Z | X_0, \mathcal{T} \in \{a, n\} \\ Y_t^d &\perp\!\!\!\perp Z | X_0, \mathcal{T} = c \\ \mathcal{T} &\perp\!\!\!\perp Z | X_0. \end{aligned}$$

Assumption 2 requires that conditional on the observed baseline characteristics, the instrumental variable Z affects neither the histories of characteristics (possibly including intermediate potential outcomes) nor of responses of the always and never takers up to one period prior to the outcome period considered. In the two outcome periods case for instance, only X_1, R_1 are restricted in this way, but not R_2 . Furthermore, note that for the always takers, the exclusion restriction only refers to the potential outcome under treatment, because $(Y_t, \mathbf{X}_{t-1}, \mathbf{R}_{t-1}) = (Y_t^1, \mathbf{X}_{t-1}^1, \mathbf{R}_{t-1}^1)$ for $\mathcal{T} = a$, while $(Y_t^0, \mathbf{X}_{t-1}^0, \mathbf{R}_{t-1}^0)$ is not restricted. An analogous statement holds for the never takers.

Assuming that Z does not affect the response behavior of always and never takers may appear reasonable in double-blind randomized medical trials, where individuals are not even aware of their treatment assignment. In non-blinded trials, this assumption seems generally less innocuous. Consider e.g. a non-blinded randomized drug-trial, where a never taker does not take the new drug irrespective of being assigned to treatment or control. Under assignment to treatment she actively decides to not comply with the protocol, whereas she would comply when being assigned to the control group. It is conceivable that the decision to not comply might affect response behavior. In other cases it may however be less of a problem. E.g. assume the randomization of a school voucher (for tuition fees) where the outcome of interest is some test score in the final grade and non-response is characterized by dropping out from school. Here, it appears more reasonable that mere voucher assignment does not affect the drop out decision of never takers, who would not use the school voucher anyway.

The stronger exclusion restriction is only required for the always and never takers, *not* for the compliers. Concerning the latter, only the standard exclusion restriction $Y_t^d \perp\!\!\!\perp Z | X_0, \mathcal{T} = c$ is imposed (see second line of Assumption 2) such that non-response may be arbitrarily related to and thus, affected by the instrument. This may happen either directly, e.g. when Z is treatment assignment and the notification of having been assigned to the treatment or control group itself changes the response behavior, or indirectly via treatment choice, e.g. due to the side effects or adverse events of a drug treatment which influences attrition.

Assumption 3: Monotonicity and support restrictions

$$\begin{aligned} \text{Existence of compliers:} & \quad \Pr(\mathcal{T} = c) > 0 \\ \text{Monotonicity:} & \quad \Pr(\mathcal{T} = d) = 0 \\ \text{Common support:} & \quad 0 < \Pr(Z = 1 | X_0) < 1. \end{aligned}$$

Assumption 3 invokes weak monotonicity, i.e. the existence of compliers and the non-existence of defiers (or vice versa). For nonparametric identification, common support in the baseline characteristics X_0 across the populations receiving and not receiving the instrument must also hold. This is e.g. satisfied in randomized experiments, where $\Pr(Z = 1 | X_0)$ is often a constant.

Theorem 1 shows the identification of the mean potential outcomes of the compliers. For ease of exposition, only two outcome periods are considered here, i.e. $t \in \{1, 2\}$, while the general result for more than two periods is provided in the online appendix. For a concise exposition of the results, we define the following conditional probabilities:

$$\begin{aligned} \pi &= \Pr(Z = 1 | X_0) \\ P_t &= \Pr(Z = 1 | X_0, \mathbf{X}_t, \mathbf{R}_t = 1, D = 1) \\ P'_t &= \Pr(Z = 1 | X_0, \mathbf{X}_t, \mathbf{R}_{t+1} = 1, D = 1) \\ \Xi_t &= \Pr(R_{t+1} = 1 | X_0, \mathbf{X}_t, \mathbf{R}_t = 1, D = 1) \\ \Xi_{t,Z=z} &= \Pr(R_{t+1} = 1 | X_0, \mathbf{X}_t, \mathbf{R}_t = 1, D = 1, Z = z). \end{aligned}$$

Identification is based on a weighting representation in which four conditional probabilities enter multiplicatively: The probability that Z takes the value one, conditional on three different sets of regressors, and a time-varying conditional response probability. For identification, Ξ_t has to be larger than zero, i.e. for each value of the covariates (X_0, \mathbf{X}_t) , the probability of attrition must not be one. Then, the treatment effect on the compliers is identified as $E[Y_t^1 | \mathcal{T} = c] - E[Y_t^0 | \mathcal{T} = c]$.

The intuition underlying Theorem 1 is as follows. By the independence of Z and \mathcal{T} given X_0 stated in Assumption 2, the proportions of compliers, always takers, and never takers in groups defined by D and Z are identified. By Assumption 1, the first period potential outcomes are independent of first period response conditional on X_0 , Z , and \mathcal{T} , and in the second period, independence of Y_2 and R_2 holds by additionally conditioning on X_1 and R_1 . Together with the exclusion restrictions on the compliers' potential outcomes as well as the potential outcomes and pre-period responses (only relevant for the second period) of always and never takers postulated in Assumption 2, this ultimately allows isolating the mean potential outcomes of compliers in the mixed groups with $(Z = 1, D = 1)$ and $(Z = 0, D = 0)$, so that the CACE is identified. Finally, it is worth noting that if there was no attrition, the CACE based on the expressions in Theorem 1 would simplify to equation (11) in Frölich (2007), which provides a representation of the CACE based on inverse probability weighting in the absence of the missing outcomes problem.

Theorem 1 *Under Assumptions 1, 2 and 3, the potential outcomes in periods $t \in \{1, 2\}$ are identified as*

$$\begin{aligned} E [Y_{t=1}^1 | \mathcal{T} = c] &= E \left[Y_{t=1} R_{t=1} \frac{D}{\pi} \frac{Z - \pi}{1 - \pi} \frac{1}{\Xi_0} \frac{P_0 - \pi}{P'_0 - \pi} \right] \times \frac{1}{E \left[\frac{D}{\pi} \frac{Z - \pi}{1 - \pi} \right]} \\ E [Y_{t=2}^1 | \mathcal{T} = c] &= E \left[Y_{t=2} R_{t=1} R_{t=2} \frac{D}{\pi} \frac{Z - \pi}{1 - \pi} \frac{1}{\Xi_0 \Xi_1} \frac{P_0 - \pi}{P'_0 - \pi} \frac{P_1 - \pi}{P'_1 - \pi} \right] \times \frac{1}{E \left[\frac{D}{\pi} \frac{Z - \pi}{1 - \pi} \right]}. \end{aligned} \quad (2)$$

An equivalent expression for $E [Y_t^0 | \mathcal{T} = c]$ is obtained by replacing D with $1 - D$ and $D = 1$ with $D = 0$ everywhere.

2.3 Non-response under the missing at random assumption

In this section we consider an alternative identification approach, where the response process is assumed to be ignorable conditional on observed characteristics, which corresponds to a type of MAR assumption. I.e., we do no longer permit that the unobserved type \mathcal{T} is related to response behavior. This implies that only unobservables that are not related to the potential outcomes are allowed to affect attrition. Again, past values of Y may trigger non-response in the current period, but neither present nor future values of Y . As stated in Assumption 1', response behavior might depend on all past values of X , which itself could be endogenous, i.e. causally affected by the treatment.

Assumption 1': Predetermined non-response

$$Y_{t+s} \perp\!\!\!\perp R_t | X_0, \underline{X}_{t-1}, \underline{R}_{t-1}, Z, D \quad \text{for } s \geq 0. \quad (3)$$

The key difference between Assumption 1' and Assumption 1 is that the latter permits the response behavior to depend on the latent type \mathcal{T} , while the former does not. Still, Assumption 1' allows response to be a function of the received treatment, which is a relevant scenario e.g. if the treatment leads to dissatisfaction and reduces the willingness to provide outcome data. On the other hand, one can think of many frameworks where it is not the treatment receipt alone that determines response behavior but rather the unobserved type \mathcal{T} of an individual, as permitted in Assumption 1. Consider e.g. an educational intervention as analyzed in Angrist, Lang, and Oreopoulos (2009) where college students are randomly provided with services and/or financial incentives to obtain better grades. In this context, never takers who do not comply when offered a treatment might have a higher probability to drop out due to a lower commitment to this particular college or to higher education in general. Assumption 1' therefore appears to be more restrictive than Assumption 1 in many empirical applications.

On the other hand, since we need no longer condition on the latent type, the restrictions on the instrument can be relaxed somewhat. The following Assumption 2' is thus a little weaker than Assumption 2 because exclusion restrictions of the instrument on the response behavior do not have to be imposed for any type. This may be of practical relevance in randomized trials e.g. if those always takers who were not randomized into the treatment ($Z = 0$) are less inclined to respond than those with $Z = 1$ due to their discontent about having to organize the treatment receipt through alternative means. In this case, Assumption 2 is violated while Assumption 2' may still hold.

Assumption 2': Exclusion restriction: For $d \in \{0, 1\}$

$$(Y_t^d, \mathcal{T}) \perp\!\!\!\perp Z | X_0.$$

Theorem 2 gives the identification results for the compliers under MAR for the case of two outcome periods, while the general result for more outcome periods is provided in the appendix.

Theorem 2 *Under Assumptions 1', 2' and 3, the potential outcomes in periods $t \in \{1, 2\}$ are*

identified as

$$\begin{aligned}
E[Y_{t=1}^1 | \mathcal{T} = c] &= E \left[\frac{Y_{t=1} R_{t=1} D Z}{\pi \Xi_{0,Z=1}} - \frac{Y_{t=1} R_{t=1} D (1-Z)}{(1-\pi) \Xi_{0,Z=0}} \right] \times \frac{1}{E \left[\frac{D}{\pi} \frac{Z-\pi}{1-\pi} \right]}, \\
E[Y_{t=2}^1 | \mathcal{T} = c] &= E \left[\frac{Y_{t=2} R_{t=1} R_{t=2} D Z}{\pi \Xi_{0,Z=1} \Xi_{1,Z=1}} - \frac{Y_{t=2} R_{t=1} R_{t=2} D (1-Z)}{(1-\pi) \Xi_{0,Z=0} \Xi_{1,Z=0}} \right] \times \frac{1}{E \left[\frac{D}{\pi} \frac{Z-\pi}{1-\pi} \right]}. \quad (4)
\end{aligned}$$

The expression for $E[Y_t^0 | \mathcal{T} = c]$ is obtained by replacing D with $1 - D$ and $D = 1$ with $D = 0$ everywhere.

Note that the assumptions underlying Theorems 1 and 2 are partly testable. Consider first the case that attrition is zero in some outcome period (e.g. zero attrition in the first follow-up period). Our setup then collapses to the standard LATE assumptions, for which tests have been proposed by Huber and Mellace (2013) and Kitagawa (2013). Similar tests could be derived for the case with attrition. By straightforward modifications of Theorems 1 and 2 the distribution functions of the potential outcomes among compliers are identified and therefore, also the density functions. As in Kitagawa (2013), a testable implication is that the estimated potential outcome densities of compliers must not be significantly negative at any point in the outcome support, because this would indicate the failure of our identifying assumptions. As a further possibility to validate the MAR assumptions underlying Theorem 2, one may consider the approach of Hirano, Imbens, Ridder, and Rubin (2001) for testing MAR models in the presence of a refreshment sample. We leave the detailed derivations and analyses of such tests for future research.

3 Finite sample properties

To illustrate the behavior of the proposed estimators in finite samples we examine a small simulation study in this section. We consider the following data generating process (DGP) with, for the sake of simplicity, parsimonious specifications of the instrument, treatment, covariate, response, and outcome equations that nevertheless give an idea about which forms of attrition can be controlled for based on our identification results:

$$\begin{aligned}
X_0 &\sim \text{uniform}(0, 1), & Z &= I\{0.25X_0 + W > 0\}, & D &= I\{\alpha Z - 0.25X_0 + U_0 > 0.5\}, \\
Y_1 &= 0.5X_0 + 0.5D + \kappa DU_0 + U_1, \\
X_1 &= 0.5Y_1 + 0.5Q, \\
R_1 &= I\{0.25X_0 + 0.25D + \beta Z + \gamma I\{0.5 - \alpha < U_0 - 0.25X_0 \leq 0.5\} + \delta Y_1 + V > 0\}, \\
Y_2 &= 0.5X_0 + X_1 + D + \kappa DU_0 + U_2, \\
R_2 &= R_1 I\{0.25X_0 + 0.25X_1 + 0.25D + \beta Z + \gamma I\{0.5 - \alpha < U_0 - 0.25X_0 \leq 0.5\} - \delta Y_2 + \epsilon > 0\},
\end{aligned}$$

each of $Q, V, W, \epsilon \sim N(0, 1)$, independent of each other,

$$\text{and } \begin{pmatrix} U_0 \\ U_1 \\ U_2 \end{pmatrix} \sim N(\mu, \sigma), \text{ where } \mu = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \text{ and } \sigma = \begin{pmatrix} 1 & 0.5 & 0.5 \\ 0.5 & 1 & 0.5 \\ 0.5 & 0.5 & 1 \end{pmatrix}.$$

$\alpha, \beta, \gamma, \delta, \kappa$ are parameters of the DGP that will be varied later. $I\{\cdot\}$ denotes the indicator function which is one if its argument is true and zero otherwise. Q, V, W, ϵ are random nuisance variables that are standard normal with zero correlation. U_0, U_1, U_2 are unobserved terms in the treatment and outcome equations in various periods. Correlation among these variables causes the endogeneity problem we have to deal with: Endogeneity is caused by the fact that U_0 affects the treatment D and is also associated with the outcomes Y_1, Y_2 through its correlation with U_1 and U_2 . The response indicators R_1, R_2 are equal to one if the outcome is observed in the respective period. Attrition is modeled as an absorbing state, i.e., R_2 is necessarily zero if $R_1 = 0$. X_0, X_1 are observed covariates. The uniformly distributed X_0 confounds the instrument because of its impact on Z and Y_1 and Y_2 . Therefore, the instrument is only conditionally valid given X_0 . The latter also affects response in both periods, thus causing attrition bias if not controlled for. Similarly, X_1 jointly influences R_2 and Y_2 , creating further bias in the second period. Note that X_1 is a function of Y_1 , which incorporates the idea that previous outcomes or functions thereof might be used to model attrition in the current period.

Our set up contains several tuning parameters: $\alpha, \beta, \gamma, \delta, \kappa$. In the treatment equation, α determines the strength of the instrument and thus the share of compliers. The compliers are those individuals with values of U_0 and X_0 such that $0.5 - \alpha < U_0 - 0.25X_0 \leq 0.5$. The larger α , the more individuals react to a change in the instrument by switching their treatment status. We consider two values of α : $\alpha=0.68$ and $\alpha=1.35$, resulting in complier shares of roughly 25% and 50% under effect homogeneity, respectively. β in the response equations gauges the effect of the instrument on R_1 and R_2 . If $\beta \neq 0$, the exclusion restriction of Z on response as postulated in Assumption

2 is violated and estimators based on Theorem 1 are inconsistent. γ defines the extent to which the compliers' response behavior differs from the remainder of the population (i.e., the never and always takers). To see this, remember that $I\{0.5 - \alpha < U_0 - 0.25X_0 \leq 0.5\}$ is an indicator for being a complier. For $\gamma \neq 0$, Assumption 1' is violated because response then depends on the latent types. In this case, estimators based on Theorem 2 are inconsistent. δ determines whether response is related to the outcomes of the current period, as for instance considered in Xie and Qian (2012). I.e., if $\delta \neq 0$, then R_t depends on Y_t such that neither Assumption 1 nor Assumption 1' are satisfied. Hence, estimators based on Theorems 1 or 2 are all inconsistent. Finally, κ determines whether the treatment effects are homogeneous or heterogeneous as a function of the unobservables U_0 . For $\kappa = 0$, the treatment effects are homogeneous, i.e. identical for everyone. In this case, the treatment effect is 0.5 for everyone in the first period and 1.25 in the second period. (The effect of 1.25 consists of the direct effect of D on Y_2 , which is 1.00, and the indirect effect of $0.5 \cdot 0.5 = 0.25$ running through X_1 .) For $\kappa \neq 0$, the treatment effects differ depending on the values of U_0 . Therefore, the CACE differs from the average effect in the total population because of different distributions of U_0 .

We simulate the DGP 1000 times with a sample size of 5000 observations, which is representative for many recently conducted field experiments in social sciences, see for instance Angrist, Bettinger, and Kremer (2006) and Bertrand and Mullainathan (2004). (The separate online appendix also examines other sample sizes.) We investigate the performance of the following estimators: (i) naive estimation based on mean differences in observed treated and non-treated outcomes that ignores both treatment endogeneity and attrition, (ii) CACE estimation based on equation (11) in Tan (2006) or equation (12) in Frölich (2007) that controls for endogeneity, but ignores attrition (denoted by $\hat{\omega}$), (iii) CACE estimation using expression (2) of Theorem 1 (denoted by $\hat{\theta}$), and (iv) CACE estimation using expression (4) of Theorem 2 (denoted by $\hat{\phi}$). The propensity scores in $\hat{\omega}$, $\hat{\theta}$, and $\hat{\phi}$ are estimated by local constant kernel regression (with Gaussian kernel). The bandwidths were chosen according to the nearest-neighbor-based default smoothing parameter in the R-package 'locfit', which was 0.7. (The results were similar when using a different kernel function such as the Epanechnikov kernel and/or when using other bandwidth values such as 0.6 and 0.8. However, values smaller than 0.6 considerably increased the variance of $\hat{\theta}$, whereas the estimates and standard errors were fairly robust for larger bandwidth values, e.g. 1.0 or larger.)

We also consider trimmed versions of $\hat{\theta}$ and $\hat{\phi}$ in order to prevent denominators from being

Table 1: Simulation 1 - treatment endogeneity and conditional LI

	Homogeneous effects											
	$\alpha=0.68, \beta=0, \gamma=0.5, \delta=0, \kappa=0$						$\alpha=1.35, \beta=0, \gamma=0.5, \delta=0, \kappa=0$					
	time period 1			time period 2			time period 1			time period 2		
	bias	stddev	rmse	bias	stddev	rmse	bias	stddev	rmse	bias	stddev	rmse
naïve	0.70	0.03	0.70	0.92	0.06	0.92	0.56	0.03	0.56	0.69	0.06	0.69
$\hat{\omega}$	0.05	0.11	0.13	0.12	0.16	0.20	0.03	0.06	0.07	0.07	0.08	0.11
$\hat{\theta}$	0.01	0.12	0.12	-0.03	0.17	0.17	0.01	0.06	0.06	-0.01	0.09	0.09
$\hat{\theta}_{\text{trim}}(0.15)$	0.01	0.12	0.12	-0.03	0.17	0.17	0.01	0.06	0.06	-0.01	0.09	0.09
$\hat{\theta}_{\text{trim}}(0.01)$	0.01	0.12	0.12	-0.03	0.17	0.17	0.01	0.06	0.06	-0.01	0.09	0.09
$\hat{\phi}$	-0.14	0.14	0.20	-0.41	0.23	0.47	-0.09	0.07	0.12	-0.25	0.11	0.27
$\hat{\phi}_{\text{trim}}(0.15)$	-0.14	0.14	0.20	-0.41	0.23	0.47	-0.09	0.07	0.12	-0.25	0.11	0.27
$\hat{\phi}_{\text{trim}}(0.01)$	-0.14	0.14	0.20	-0.41	0.23	0.47	-0.09	0.07	0.12	-0.25	0.11	0.27
MAR G-comp.	0.71	0.03	0.71	0.94	0.06	0.94	0.56	0.03	0.56	0.70	0.06	0.70
Heckman	0.36	1.59	1.63	0.03	1.61	1.61	0.33	3.30	3.32	-0.22	2.63	2.64
true CACE		0.50			1.25			0.50			1.25	
mean response		0.63			0.41			0.69			0.49	

	Heterogeneous effects											
	$\alpha=0.68, \beta=0, \gamma=0.5, \delta=0, \kappa=0.5$						$\alpha=1.35, \beta=0, \gamma=0.5, \delta=0, \kappa=0.5$					
	time period 1			time period 2			time period 1			time period 2		
	bias	stddev	rmse	bias	stddev	rmse	bias	stddev	rmse	bias	stddev	rmse
naïve	0.99	0.04	0.99	1.34	0.06	1.34	0.83	0.04	0.83	1.10	0.06	1.10
$\hat{\omega}$	0.05	0.13	0.14	0.12	0.18	0.22	0.03	0.07	0.08	0.08	0.10	0.13
$\hat{\theta}$	0.01	0.13	0.13	-0.03	0.20	0.20	0.01	0.07	0.07	-0.01	0.10	0.10
$\hat{\theta}_{\text{trim}}(0.15)$	0.01	0.13	0.13	-0.03	0.20	0.20	0.01	0.07	0.07	-0.01	0.10	0.10
$\hat{\theta}_{\text{trim}}(0.01)$	0.01	0.13	0.13	-0.03	0.20	0.20	0.01	0.07	0.07	-0.01	0.10	0.10
$\hat{\phi}$	-0.20	0.16	0.25	-0.52	0.26	0.58	-0.14	0.08	0.16	-0.35	0.13	0.37
$\hat{\phi}_{\text{trim}}(0.15)$	-0.20	0.16	0.25	-0.52	0.26	0.58	-0.14	0.08	0.16	-0.35	0.13	0.37
$\hat{\phi}_{\text{trim}}(0.01)$	-0.20	0.16	0.25	-0.52	0.26	0.58	-0.14	0.08	0.16	-0.35	0.13	0.37
MAR G-comp.	1.00	0.04	1.00	1.36	0.06	1.37	0.84	0.04	0.84	1.10	0.06	1.10
Heckman	0.63	2.27	2.36	0.25	2.07	2.09	0.41	5.04	5.06	-0.18	2.85	2.86
true CACE		0.64			1.45			0.48			1.22	
mean response		0.63			0.43			0.69			0.50	

Note: Results are based on 1000 simulations and 5000 observations.

close to zero, which may imply arbitrarily large weights for some observations. Propensity score trimming is discussed e.g. in Frölich (2004), Heckman, Ichimura, and Todd (1997), Dehejia and Wahba (1999), Busso, DiNardo, and McCrary (2009), and Crump, Hotz, Imbens, and Mitnik (2009). Yet, a trimming rule that is optimal in the sense that it minimizes the mean square error of the estimator does not appear to be available in the literature. Here, we follow Huber, Lechner, and Wunsch (2013) and discard observations whose relative weights within subgroups defined by Z and D exceed a particular threshold. As trimming thresholds we consider relative weights of 15 and 1%, resulting in the trimmed estimators $\hat{\theta}_{\text{trim}}(0.15)$, $\hat{\phi}_{\text{trim}}(0.15)$, $\hat{\theta}_{\text{trim}}(0.01)$, $\hat{\phi}_{\text{trim}}(0.01)$. The appendix provides additional results for further trimming levels (10, 5, and 2%).

We also consider an estimator that controls for attrition under the assumption of MAR but ignores treatment endogeneity due to U_0, U_1, U_2 , while controlling for confounding related to X_0 . To be specific, we use the MLE-based G-computation procedure of Robins (1986), in which the outcomes and response processes are modeled parametrically by linear and logit specifications, respectively. The appendix also provides the results for estimation based on targeted MLE, see van der Laan and Rubin (2006), inverse probability weighting (see e.g. Horvitz and Thompson (1952) and Hirano, Imbens, and Ridder (2003)), and augmented IPW (AIPW) (as in Robins, Rotnitzky, and Zhao (1995) and Scharfstein, Rotnitzky, and Robins (1999)) which yield very similar results. Finally, parametric Heckman (1976) MLE estimation of sample selection models assuming jointly normally distributed unobserved terms in the response and the outcome equations is also considered. The latter estimator controls for X_0, D , and Z in the estimation of response and can therefore account for attrition related to unobservables if R_1 and R_2 are functions of Z and if Z does not have a direct effect on the outcomes conditional on X_0 and D . However, it does not allow for treatment endogeneity related to U_0, U_1, U_2 and additionally presumes treatment effects to be homogeneous.

Table 1 provides the bias, standard deviation and root mean squared error (rmse) of the various estimators in periods 1 and 2 under treatment endogeneity and conditional LI with $\gamma = 0.5$, and β, δ equal to zero. $\hat{\theta}$, which is consistent in this scenario, performs very well in terms of bias and rmse irrespective of the period, share of compliers and effect homogeneity or heterogeneity. In contrast, the naive approach, the MAR-based G-computation procedure not controlling for treatment endogeneity, and the Heckman estimator are severely biased in any specification. Also $\hat{\phi}$ (and its trimmed versions) and $\hat{\omega}$ are prone to non-negligible bias, even though the latter performs

Table 2: Simulation 2 - treatment endogeneity and MAR (with the instrument affecting response)

	Homogeneous effects											
	$\alpha=0.68, \beta=0.5, \gamma=0, \delta=0, \kappa=0$						$\alpha=1.35, \beta=0.5, \gamma=0, \delta=0, \kappa=0$					
	time period 1			time period 2			time period 1			time period 2		
	bias	stddev	rmse	bias	stddev	rmse	bias	stddev	rmse	bias	stddev	rmse
naive	0.75	0.03	0.75	1.08	0.06	1.09	0.64	0.03	0.64	0.93	0.06	0.93
$\hat{\omega}$	0.21	0.09	0.23	0.47	0.12	0.48	0.12	0.05	0.13	0.28	0.07	0.29
$\hat{\theta}$	-0.37	41.34	41.34	307.31	9507.59	9512.56	-0.05	0.08	0.10	-0.40	4.80	4.82
$\hat{\theta}_{\text{trim}}(0.15)$	-0.93	5.08	5.17	1.76	32.97	33.02	-0.05	0.08	0.10	-0.49	0.64	0.81
$\hat{\theta}_{\text{trim}}(0.01)$	-0.86	4.30	4.38	2.10	41.21	41.26	-0.05	0.08	0.10	-0.43	0.84	0.95
$\hat{\phi}$	0.02	0.14	0.14	-0.09	0.23	0.25	0.01	0.07	0.07	-0.05	0.12	0.13
$\hat{\phi}_{\text{trim}}(0.15)$	0.02	0.14	0.14	-0.09	0.23	0.25	0.01	0.07	0.07	-0.05	0.12	0.13
$\hat{\phi}_{\text{trim}}(0.01)$	0.02	0.14	0.14	-0.09	0.23	0.25	0.01	0.07	0.07	-0.05	0.12	0.13
MAR G-comp.	0.76	0.03	0.77	1.11	0.05	1.11	0.65	0.03	0.65	0.95	0.06	0.95
Heckman	0.86	0.04	0.86	1.30	0.06	1.30	0.83	0.04	0.83	1.26	0.09	1.26
true CACE	0.50			1.25			0.50			1.25		
mean response	0.68			0.49			0.69			0.51		

	Heterogeneous effects											
	$\alpha=0.68, \beta=0.5, \gamma=0, \delta=0, \kappa=0.5$						$\alpha=1.35, \beta=0.5, \gamma=0, \delta=0, \kappa=0.5$					
	time period 1			time period 2			time period 1			time period 2		
	bias	stddev	rmse	bias	stddev	rmse	bias	stddev	rmse	bias	stddev	rmse
naive	1.05	0.04	1.05	1.54	0.06	1.54	0.94	0.04	0.94	1.39	0.06	1.39
$\hat{\omega}$	0.29	0.10	0.31	0.65	0.13	0.67	0.18	0.06	0.19	0.43	0.09	0.44
$\hat{\theta}$	-0.28	41.34	41.34	307.47	9507.59	9512.56	0.01	0.09	0.09	-0.27	4.80	4.81
$\hat{\theta}_{\text{trim}}(0.15)$	-0.85	5.08	5.15	1.92	32.97	33.03	0.01	0.09	0.09	-0.36	0.65	0.74
$\hat{\theta}_{\text{trim}}(0.01)$	-0.77	4.23	4.37	2.26	41.21	41.27	0.01	0.09	0.09	-0.30	0.84	0.89
$\hat{\phi}$	0.02	0.15	0.15	-0.08	0.25	0.26	0.01	0.08	0.08	-0.04	0.13	0.14
$\hat{\phi}_{\text{trim}}(0.15)$	0.02	0.15	0.15	-0.08	0.25	0.26	0.01	0.08	0.08	-0.04	0.13	0.14
$\hat{\phi}_{\text{trim}}(0.01)$	0.02	0.15	0.15	-0.08	0.25	0.26	0.01	0.08	0.08	-0.04	0.13	0.14
MAR G-comp.	1.07	0.04	1.07	1.57	0.06	1.58	0.95	0.04	0.95	1.41	0.06	1.41
Heckman	1.21	0.04	1.21	1.88	0.07	1.88	1.19	0.04	1.19	2.02	0.44	2.06
true CACE	0.64			1.45			0.48			1.22		
mean response	0.68			0.49			0.69			0.51		

Note: Results are based on 1000 simulations and 5000 observations.

comparably well in the first time period. Note that trimming does neither affect $\hat{\theta}$, nor $\hat{\phi}$, implying that large relative weights do not occur.

In the second simulation (Table 2), $\gamma = 0$ such that the assumptions underlying $\hat{\phi}$ hold. At the same time $\beta = 0.5$, implying a direct effect of the instrument on the response process and a violation of Assumption 2 required for the consistency of $\hat{\theta}$. Hence, estimators based on Theorem 2 are consistent, whereas the assumptions for Theorem 1 are not met. As expected, $\hat{\phi}$ now dominates any other estimator with respect to bias and low rmse and is unchanged by trimming. The naive approach, $\hat{\theta}$, G-computation, the Heckman estimator, and (to a lesser extent) $\hat{\omega}$ are substantially biased in most cases. $\hat{\theta}$ performs particularly poorly under the smaller complier share ($\alpha = 0.68$) due to a large increase of the variance. Yet, already moderate trimming using the 15% threshold ($\hat{\theta}_{\text{trim}}(0.15)$) reduces the variance (and the rmse) considerably, even though it remains at comparably high levels. More trimming further decreases the rmse in the first period, but increases it in the second one. In the latter case, the rmse is relatively stable for 15% and 10%, but grows more strongly for 2% and 1%.

In the third simulation (Table 3) we consider a scenario where all estimators are inconsistent: γ is set to zero, while $\beta = 0.5$ and $\delta = 0.25$, implying that the instrument directly affects non-response, which in addition is also related to the outcomes of the current period. $\hat{\theta}$ and $\hat{\phi}$ are biased because they ignore attrition related to contemporaneous outcomes, while G-computation ignores both treatment endogeneity and attrition related to contemporaneous outcomes, and the Heckman estimator does not account for treatment endogeneity. Trimming again reduces the variance of $\hat{\theta}$ in several cases, but smaller threshold values tend to increase the rmse relative to larger thresholds when $\alpha=0.68$. All in all, no method performs convincingly in this last set-up considered.

4 Application to a policy intervention in college

In this section, we apply our methods to data from the Student Achievement and Retention Project assessed in Angrist, Lang, and Oreopoulos (2009), a randomized program providing academic services and financial incentives to first year students at a Canadian campus which aimed at improving the academic performance. To this end, all students who entered in September 2005 and had a high school grade point average (GPA) lower than the upper quartile were randomly assigned either to one of three different treatments provided in the first year, namely academic support services, financial incentives, or both, or otherwise to a control group. The services

Table 3: Simulation 3 - treatment endogeneity and selection on current outcomes

	Homogeneous effects											
	$\alpha=0.68, \beta=0.5, \gamma=0, \delta=0.25, \kappa=0$						$\alpha=1.35, \beta=0.5, \gamma=0, \delta=0.25, \kappa=0$					
	time period 1			time period 2			time period 1			time period 2		
	bias	stddev	rmse	bias	stddev	rmse	bias	stddev	rmse	bias	stddev	rmse
naive	0.68	0.03	0.68	0.85	0.05	0.85	0.57	0.03	0.57	0.67	0.06	0.67
$\hat{\omega}$	0.10	0.10	0.13	0.36	0.11	0.37	0.05	0.05	0.07	0.19	0.07	0.20
$\hat{\theta}$	-1.55	5.09	5.32	1.48	572.82	572.82	-0.19	0.08	0.21	-0.68	0.47	0.83
$\hat{\theta}_{\text{trim}}(0.15)$	-1.18	4.65	4.80	-1.21	20.30	20.34	-0.19	0.08	0.21	-0.67	0.24	0.71
$\hat{\theta}_{\text{trim}}(0.01)$	-1.21	4.90	5.05	-0.45	24.37	24.37	-0.19	0.08	0.21	-0.65	0.20	0.69
$\hat{\phi}$	-0.23	0.14	0.27	-0.59	0.22	0.63	-0.13	0.07	0.15	-0.36	0.11	0.37
$\hat{\phi}_{\text{trim}}(0.15)$	-0.23	0.14	0.27	-0.59	0.22	0.63	-0.13	0.07	0.15	-0.36	0.11	0.37
$\hat{\phi}_{\text{trim}}(0.01)$	-0.23	0.14	0.27	-0.59	0.22	0.63	-0.13	0.07	0.15	-0.36	0.11	0.37
MAR G-comp.	0.70	0.03	0.70	0.89	0.05	0.89	0.58	0.03	0.58	0.70	0.05	0.70
Heckman	0.87	0.04	0.88	1.23	0.09	1.23	0.84	0.07	0.84	2.11	1.02	2.34
true CACE	0.50			1.25			0.50			1.25		
mean response	0.71			0.55			0.72			0.57		
	Heterogeneous effects											
	$\alpha=0.68, \beta=0.5, \gamma=0, \delta=0.25, \kappa=0.5$						$\alpha=1.35, \beta=0.5, \gamma=0, \delta=0.25, \kappa=0.5$					
	time period 1			time period 2			time period 1			time period 2		
	bias	stddev	rmse	bias	stddev	rmse	bias	stddev	rmse	bias	stddev	rmse
naive	1.00	0.03	1.00	1.33	0.06	1.33	0.89	0.04	0.89	1.18	0.06	1.18
$\hat{\omega}$	0.13	0.11	0.17	0.43	0.14	0.45	0.09	0.06	0.11	0.28	0.09	0.29
$\hat{\theta}$	-1.51	5.09	5.31	1.55	572.82	572.82	-0.15	0.09	0.18	-0.60	0.47	0.76
$\hat{\theta}_{\text{trim}}(0.15)$	-1.14	4.65	4.79	-1.13	20.30	20.33	-0.15	0.09	0.18	-0.59	0.25	0.64
$\hat{\theta}_{\text{trim}}(0.01)$	-1.16	4.90	5.04	-0.38	24.37	24.37	-0.15	0.09	0.18	-0.57	0.21	0.61
$\hat{\phi}$	-0.22	0.15	0.27	-0.56	0.24	0.61	-0.11	0.08	0.13	-0.29	0.12	0.31
$\hat{\phi}_{\text{trim}}(0.15)$	-0.22	0.15	0.27	-0.56	0.24	0.61	-0.11	0.08	0.13	-0.29	0.12	0.31
$\hat{\phi}_{\text{trim}}(0.01)$	-0.22	0.15	0.27	-0.56	0.24	0.61	-0.11	0.08	0.13	-0.29	0.12	0.31
MAR G-comp.	1.02	0.03	1.02	1.37	0.06	1.37	0.90	0.04	0.90	1.21	0.06	1.21
Heckman	1.30	0.04	1.30	1.95	0.27	1.97	1.67	0.57	1.77	2.65	1.27	2.94
true CACE	0.64			1.45			0.48			1.22		
mean response	0.72			0.56			0.73			0.59		

Note: Results are based on 1000 simulations and 5000 observations.

contained both access to peer advisors, i.e., trained upper-class students supposed to provide academic support, and class-specific sessions targeted at improving study habits without focusing on specific course content. The financial incentives consisted of cash payments between 1,000 and 5,000 dollars that were conditional on attaining particular GPA targets in college, where the targets were a function of the high school GPA.

While the intervention appeared to be generally ineffective for males, Angrist, Lang, and Oreopoulos (2009) found positive effects of the combined treatment (academic support and financial incentives) on the college performance of females in the first and second year. For this reason, we will only focus on the subsample of 948 female students in the subsequent discussion. As the number of observations assigned to a particular treatment arm is rather low, we aggregate the academic services and financial incentives to a binary treatment that takes the value one if any form of intervention took place and zero otherwise in order to avoid small sample problems. For the same reason, we use (parametric) probit regressions (rather than nonparametric methods) to estimate the conditional probabilities involved in the identification results, which entails semiparametric estimators of the CACE. Inference is based on the bootstrap.

Albeit treatment assignment was random, identification may be flawed by both endogeneity and attrition. The endogeneity issue stems from the fact that only 274 (or 73%) of the 374 students who were offered any treatment actually signed up for it, which gives rise to potential selection bias into treatment. Furthermore, GPA scores, one of the outcomes measuring college success, are not observed for all students. Whereas they are missing for only 56 students (or 6%) in the first year, non-response amounts to a non-negligible 169 (or 18%) in the second year. If attrition is selective so that e.g. the probability to drop out decreased in both the treatment state and unobserved ability, the treatment effect is biased due to positive selection into observed GPA scores. Angrist, Lang, and Oreopoulos (2009) use instrumental variable estimation to control for endogeneity, where the random assignment indicator serves as instrument. They, however, do not correct for attrition in the GPA outcomes, but merely base their analysis on all those observations without missing GPAs, see the note underneath Table 6 in their paper. Here, we apply the methods outlined in Sections 2.2 and 2.3 to control for both endogeneity *and* attrition.

We are interested in the effect of having signed up for any of the three treatments ($D = 1$) vs. no treatment ($D = 0$) on the GPA scores at the end of the first and second year. We estimate the CACE based on Theorem 1 to allow attrition to be related to the latent types, as compliers with

the treatment assignment may be more motivated to stay in college than the never takers, whose reluctance to take the treatment even when offered may be associated with a higher inclination to drop out of college. This motivates our higher confidence in Assumption 1 rather than the stronger Assumption 1' (which does not permit LI conditional on observables). At the same time, it seems likely that mere assignment does not affect the drop out decision of never takers, who would not take advantage of the treatment anyway. We therefore suspect Assumption 2 to be satisfied, albeit somewhat stronger than Assumption 2'. Nevertheless, we also consider estimation based on Theorem 2 imposing MAR given the observed variables and the treatment, which allows checking the sensitivity of the results to the presumed form of attrition. If one obtains similar results under both methods, this may imply that (the respective stronger assumption of) both sets of assumptions are satisfied, i.e. Assumption 1' and Assumption 2. We use both untrimmed and trimmed versions of the respective estimators. As in the simulations, trimming discards observations whose relative weights in subgroups defined by Z and D exceed a certain threshold, which is set to 10% in the application.

The data set contains a range of pre-treatment variables measuring performance and ambition as well as socioeconomic characteristics that allow us to model the response process in the first year. E.g., we observe the GPA score in high school, the fall grade of the first year, and the attempted maths and science courses, which are most likely correlated with both GPA scores in later periods and the probability to drop out. Indeed, the empirical relevance of academic performance in high school and in the first semester of college as a predictor for attrition is well documented in the literature on higher education, see e.g. Leppel (2002), Herzog (2005), and Tinto (1997). Furthermore, the data includes self-assessed measures of effort and ambition, e.g., whether the student wants to finish in four years, or strives for a higher degree than a BA. Learning habits are reflected by the information on how often a student leaves studying until the last minute. The data also comprises important characteristics reflecting the socioeconomic background, such as age, parents' education, and indicators for living at home and English mother tongue. Finally, it contains dummies for whether the student is at the first choice college and whether she completed the base line survey which may be correlated with the likelihood to be observed in later periods.

Table 4 gives the results of a probit regression of first year response on the baseline covariates X_0 and the treatment indicator D in order to estimate $\Xi_0 = \Pr(R_1 = 1|X_0, D = 1)$. The main specification (1) contains all regressors, and is used for the estimation of the CACE. Specification

Table 4: Probit coefficients and marginal effects of the model for 1st year response

	Coefficients			Marginal effects		
	(1)	(2)	(3)	(1)	(2)	(3)
Constant	5.495 (2.896)	6.196 (2.721)	1.445 (0.072)			
Treatment D	0.335 (0.265)	0.529 (0.262)	0.571 (0.184)	0.003 (0.002)	0.004 (0.003)	0.052 (0.013)
High school GPA ≤ 75.2	-0.574 (0.312)	-0.503 (0.300)		-0.008 (0.008)	-0.007 (0.007)	
High school GPA	-0.078 (0.036)	-0.076 (0.034)		-0.001 (0.001)	-0.001 (0.001)	
Fall grade	0.028 (0.003)	0.029 (0.003)		0.000 (0.000)	0.000 (0.000)	
Attempted math/science credits	0.878 (0.239)	0.970 (0.247)		0.008 (0.004)	0.009 (0.005)	
Wants more than B.A.	0.235 (0.223)			0.002 (0.003)		
Last minute learning (usual/often)	-0.054 (0.244)			-0.000 (0.002)		
Age < 20	0.311 (0.396)			0.004 (0.008)		
Father has college degree	-0.129 (0.230)			-0.001 (0.002)		
At first choice college	0.208 (0.258)			0.002 (0.002)		
Completed baseline survey	0.778 (0.269)			0.018 (0.016)		
Pseudo R^2	0.494	0.452	0.027	0.494	0.452	0.027

Note: (1)-(3) give the probit coefficients and marginal effects, respectively, when estimating $\Xi_0 = \Pr(R_1 = 1|X_0, D = 1)$

under different specifications: (1) is the main specification with all regressors, (2) contains a subset of regressors, (3) contains only D and a constant. The marginal effects are evaluated at the means of all other regressors. Standard errors are given in

brackets. The sample size is 948.

(2) presents a more parsimonious model consisting of D and pre-treatment outcomes (high school GPA, fall grade, attempted math/science credits). Finally, specification (3) only contains D and a constant as regressors. Comparing the results for the different specifications, we find that the pre-treatment outcomes and the dummy for survey completed clearly have the highest predictive power, whereas socioeconomic variables are less important.

For modeling response in the second period, we use in addition to the covariates X_0 of specification (1) three intermediate outcomes (X_1) at the end of the first year: the GPA as well as the number of credits earned in the first year and an indicator for good standing, all of which are highly correlated with response in the second year.

Table 5 provides descriptive statistics (means and standard deviations) of the variables used in our analysis for all females, as well as for subsamples with $D = 1$ and $D = 0$. The variables measured after the first or second year are only observed if $R_1 = 1$ and $R_2 = 1$, respectively. Note that treated females have on average higher pre-treatment outcomes (high school GPA and fall grade) and higher aspirations (wanting more than a B.A.) than the non-treated. This points to selectivity and motivates the use of random treatment assignment Z as an instrument for actual treatment take-up D .

Table 5: Descriptive statistics

Regressor	total sample (948 obs.)		$D = 1$ (274 obs.)		$D = 0$ (674 obs.)	
	mean	std. dev	mean	std. dev	mean	std. dev
High school GPA (multi-valued)	78.88	4.29	79.10	4.30	78.80	4.28
Fall grade (multi-valued)	53.69	25.71	58.55	23.04	51.71	26.48
Attempted math/science credits (multi-valued)	1.00	1.16	1.05	1.19	0.97	1.15
Wants more than B.A. (binary)	0.52	0.50	0.58	0.49	0.49	0.50
Last minute learning (binary)	0.28	0.45	0.30	0.46	0.28	0.45
At first choice college (binary)	0.24	0.43	0.26	0.44	0.23	0.42
Age < 20 (binary)	0.97	0.17	0.99	0.12	0.97	0.18
Father has college degree (binary)	0.37	0.48	0.40	0.49	0.36	0.48
Completed baseline survey (binary)	0.90	0.30	0.95	0.23	0.89	0.32
First year response R_1 (binary)	0.94	0.24	0.98	0.15	0.93	0.26
First year GPA Y_1 (multi-valued)	1.76	0.90	1.81	0.88	1.74	0.91
First year good standing for $R_1 = 1$ (binary)	0.48	0.50	0.54	0.50	0.46	0.50
First year credits earned for $R_1 = 1$ (multi-valued)	2.36	0.93	2.47	0.94	2.32	0.92
Second year response R_2 (binary)	0.82	0.38	0.83	0.37	0.82	0.39
Second year GPA Y_2 for $R_2 = 1$ (multi-valued)	2.07	0.87	2.19	0.86	2.01	0.87

Note: Descriptive statistics for baseline covariates X_0 , response indicators R_1 and R_2 and outcomes Y_1 and Y_2 , if observed.

Table 6 presents the estimated treatment effects of the intervention on the GPA one and two

years later. The top panel provides the estimates for the full sample. The subsequent panels show the results for various subsamples defined by age and parental background. In each panel, the first line gives the CACE estimates, the second line the bootstrap standard errors, and the third line the bootstrap p-values based on the quantiles of the resampled distribution of the CACE estimates, see equation (6) in MacKinnon (2006). We provide the quantile-based p-values (rather than those based on the t-statistic) to account for the problem that in finite samples the moments of instrumental variable estimators may not exist such that t-statistics may be misleading, which might even be aggravated by attrition. The first and sixth columns labelled "Wald" show the Wald estimates, i.e. the instrumental variable estimator without any covariates. The estimates based on Theorem 1 are denoted by $\hat{\theta}$ and $\hat{\theta}_{\text{trim}}$, where the latter represents the trimmed version. The estimates based on Theorem 2 are denoted by $\hat{\phi}$ and $\hat{\phi}_{\text{trim}}$. We find that large weights rarely occur such that the trimmed and untrimmed point estimates are always very similar, if not the same. Note, however, that trimming reduces the standard errors of the estimates based on Theorem 1 by disciplining outliers in the bootstrap samples.

Both $\hat{\theta}$ and $\hat{\theta}_{\text{trim}}$ are nevertheless less precise than $\hat{\phi}$ and $\hat{\phi}_{\text{trim}}$. We would generally (and specifically in moderate samples) expect this to be the case at least if both theorems are (closely) satisfied, because Theorem 1 contains more conditional probabilities to be estimated, e.g. $P'_0 - \pi$ and $P'_1 - \pi$ in the denominator, which may potentially decrease precision in small samples. In particular, if the latter differences are small (which likely occurs if Z only weakly shifts D so that few compliers exist) the variance might be large. Furthermore, in the current application, $\hat{\phi}$ and $\hat{\phi}_{\text{trim}}$ appear to rest on stronger assumptions than $\hat{\theta}$ and $\hat{\theta}_{\text{trim}}$, which again suggests lower standard errors of the former: Whereas we argued in Section 2 that Assumption 2' is generally weaker than Assumption 2, they are, however, very similar in the application at hand. This is because Assumption 2 only restricts the response process in time period 1, where we have in fact very little non-response. (Non-response is larger in time period 2, but this does not enter Assumption 2.) On the other hand, Assumption 1' is clearly much stronger than Assumption 1. The former imposes independence within each stratum defined by Z and D (and other pre-determined observables), whereas the latter additionally requires conditioning on the (unobserved) type. Therefore, estimators based on Assumption 1' exploit more restrictions and can (figuratively speaking) use coarser strata with more information than methods relying on Assumption 1, which have to operate within finer strata additionally defined upon the type. Therefore, $\hat{\phi}$ and $\hat{\phi}_{\text{trim}}$ can

exploit more information.

Examining first the estimates for the whole population, we do not find any significant effects in the first year. In contrast, the simple Wald estimates for the second year are significant (at the 5% level) and suggest that the GPA of compliers increases by 0.164 points. However, when using the attrition corrected estimators, the effect shrinks considerably to 0.077 or 0.071, respectively, and becomes insignificant. Therefore, our results suggest that attrition, if ignored, may lead to an overestimation of the effects in education experiments.

In the remainder of Table 6, we investigate effect heterogeneity for subsamples stratified by age, prior academic achievement and parental background. E.g., we separately consider students in the lower and the upper half of the high school GPA distribution (median: 78.5 points) to see whether high or low achievers particularly benefit from the intervention. Indeed, the Wald estimate for the second year GPA of low achievers amounts to 0.225 points, indicating that the less capable students benefit most when taking advantage of the services and incentives. However, after controlling for attrition, the effect becomes much smaller and insignificant, irrespective of trimming. When we split the sample by age groups (17 & 18 years versus older than 18), we also cannot draw reliable conclusions as the estimates are generally rather noisy.

Finally, we examine whether the effects differ by parents' education, which might be regarded as a proxy for family background. Interestingly, the second year Wald estimate in the subsample with mothers that have a college degree is negative and large. When controlling for attrition, the estimate shrinks in magnitude (in the case of $\hat{\theta}$, $\hat{\theta}_{\text{trim}}$ quite considerably) and becomes even less significant. In contrast, for those students whose mother has no degree, the Wald estimate is significantly positive (at the 5% level) in both periods. Furthermore, correcting for attrition does *not* substantially reduce the estimate in the second year, even though the precision decreases. The estimates $\hat{\phi}$ and $\hat{\phi}_{\text{trim}}$ remain significant at the 5% level. A similar pattern appears when stratifying on the father's degree status. While the Wald estimate in the second year is insignificant in the subpopulation with fathers having a degree, it is large and significant in the subsample without college degree. Furthermore, the effect is almost the same when using $\hat{\theta}$, albeit less precisely estimated, and $\hat{\theta}_{\text{trim}}$, $\hat{\phi}$, $\hat{\phi}_{\text{trim}}$ are significant at the 10% and 5% levels, respectively.

In summary, our findings suggest that the empirical evidence about the effectiveness of the intervention considered by Angrist, Lang, and Oreopoulos (2009) is much weaker once attrition is acknowledged. Nevertheless, female students with a less favorable family background seem to gain

Table 6: Effectiveness of the school intervention on GPA: all females and subsamples

	1st year effect of intervention on GPA					2nd year effect of intervention on GPA				
	Wald estimate (no covariates)	$\hat{\theta}$ (Theorem 1)	$\hat{\theta}_{\text{trim}}$	$\hat{\phi}$ (Theorem 2)	$\hat{\phi}_{\text{trim}}$	Wald estimate (no covariates)	$\hat{\theta}$ (Theorem 1)	$\hat{\theta}_{\text{trim}}$	$\hat{\phi}$ (Theorem 2)	$\hat{\phi}_{\text{trim}}$
<i>Full sample: all females (948 obs.)</i>										
effect	0.074	0.022	0.022	-0.047	-0.047	0.164	0.077	0.077	0.071	0.071
s.e.	0.079	0.291	0.129	0.075	0.076	0.083	8.694	0.214	0.090	0.093
p-val	0.399	0.709	0.679	0.517	0.577	0.040	0.817	0.770	0.419	0.340
<i>Subsample: high school GPA ≤ 78.5 (467 obs.)</i>										
effect	0.156	0.107	0.107	0.057	0.057	0.225	0.143	0.143	0.043	0.043
s.e.	0.110	2.340	0.274	0.112	0.112	0.128	10.937	0.378	0.130	0.131
p-val	0.149	0.502	0.456	0.595	0.565	0.069	0.451	0.382	0.698	0.669
<i>Subsample: high school GPA > 78.5 (481 obs.)</i>										
effect	0.023	-0.093	-0.093	-0.170	-0.170	0.099	0.036	-0.066	-0.001	-0.122
s.e.	0.106	36.995	0.382	0.135	0.128	0.098	24.073	0.968	0.180	0.156
p-val	0.812	0.668	0.712	0.220	0.382	0.300	0.941	0.670	0.816	0.456
<i>Subsample: 17 and 18 years old (741 obs.)</i>										
effect	0.042	-0.024	-0.024	-0.050	-0.050	0.132	0.093	0.093	0.043	0.043
s.e.	0.090	2.472	0.310	0.089	0.088	0.092	28.780	0.322	0.103	0.103
p-val	0.632	0.816	0.848	0.547	0.556	0.145	0.573	0.469	0.680	0.668
<i>Subsample: 19-23 years old (207 obs.)</i>										
effect	0.131	0.064	0.002	-0.051	-0.059	0.226	0.275	0.171	0.125	-0.021
s.e.	0.192	1530.530	0.627	0.194	0.195	0.218	733.224	0.889	0.379	0.306
p-val	0.484	0.822	0.776	0.772	0.939	0.259	0.615	0.448	0.656	0.703
<i>Subsample: mother has college degree (304 obs.)</i>										
effect	-0.178	-0.171	-0.249	-0.201	-0.201	-0.253	-0.085	-0.197	-0.211	-0.211
s.e.	0.143	2.745	0.329	0.167	0.135	0.155	64.916	0.481	0.193	0.157
p-val	0.212	0.388	0.148	0.157	0.224	0.105	0.936	0.675	0.208	0.281
<i>Subsample: mother has no college degree (644 obs.)</i>										
effect	0.197	0.135	0.135	0.080	0.080	0.383	0.345	0.345	0.259	0.259
s.e.	0.097	5.960	0.098	0.111	0.098	0.099	5.095	0.271	0.155	0.131
p-val	0.034	0.130	0.117	0.366	0.290	0.000	0.147	0.096	0.050	0.018
<i>Subsample: father has college degree (355 obs.)</i>										
effect	0.035	-0.029	-0.029	-0.038	-0.038	-0.078	0.065	0.065	-0.127	-0.127
s.e.	0.133	58.422	0.216	0.187	0.123	0.141	37.131	0.406	0.206	0.134
p-val	0.827	0.839	0.886	0.717	0.833	0.566	0.457	0.212	0.342	0.443
<i>Subsample: father has no college degree (593 obs.)</i>										
effect	0.092	0.066	0.066	-0.008	-0.008	0.317	0.307	0.307	0.250	0.250
s.e.	0.101	8.636	0.133	0.099	0.101	0.108	6.418	0.225	0.140	0.143
p-val	0.371	0.491	0.479	0.979	0.904	0.002	0.150	0.085	0.063	0.029

Note: Treatment effects of the intervention (support and/or financial services) on GPA outcomes one and two years later,

respectively. The top panel displays the results for the full sample (on the left, the effect after one year; on the right, the effect after two years). The subsequent panels show estimates for subpopulations stratified by age, parental background, and prior academic achievement. P-values are given in brackets and are based on 1999 bootstrap replications. Trimming in $\hat{\theta}_{\text{trim}}$

and $\hat{\phi}_{\text{trim}}$ is based on dropping observations that have a relative weight larger than 10%.

from the services and financial incentives.

5 Conclusions

In this paper, we proposed a novel approach for the identification and estimation of local average treatment effects in multiple outcome periods which controls for both treatment endogeneity and outcome attrition. We showed how pre-treatment information can be combined with intermediate outcomes in order to correct more plausibly for non-response bias in later periods, while an instrument was used to tackle endogenous treatment selection. Two sets of identifying assumptions were presented. The first one, which we call conditional latent ignorability, permits attrition to depend on observables and the latent treatment compliance type, which may be related to unobservables. The second one imposes randomness given observed variables only, which amounts to a dynamic missing at random assumption. The proposed methods were applied to a policy intervention aimed at increasing academic performance in college, where ignoring attrition was found to lead to upwardly biased estimates.

References

- ABADIE, A. (2003): “Semiparametric instrumental Variable estimation of treatment response models,” *Journal of Econometrics*, 113, 231–263.
- ABOWD, J., B. CREPON, AND F. KRAMARZ (2001): “Moment Estimation With Attrition: An Application to Economic Models,” *Journal of the American Statistical Association*, 96, 1223–1230.
- ANGRIST, J., E. BETTINGER, AND M. KREMER (2006): “Long-Term Educational Consequences of Secondary School Vouchers: Evidence from Administrative Records in Colombia,” *American Economic Review*, 96(3), 847–862.
- ANGRIST, J., G. IMBENS, AND D. RUBIN (1996): “Identification of Causal Effects using Instrumental Variables,” *Journal of American Statistical Association*, 91, 444–472 (with discussion).
- ANGRIST, J., D. LANG, AND P. OREOPOULOS (2009): “Incentives and Services for College Achievement: Evidence from a Randomized Trial,” *American Economic Journal: Applied Economics*, 1, 136–163.
- BARNARD, J., C. FRANGAKIS, J. HILL, AND D. RUBIN (2003): “A Principal Stratification Approach to Broken Randomized Experiments: A Case Study of School Choice Vouchers in New York City,” *Journal of the American Statistical Association*, 98, 299–323.

- BERTRAND, M., AND S. MULLAINATHAN (2004): “Are Emily and Greg More Employable than Lakisha and Jamal? A Field Experiment on Labor Market Discrimination,” *The American Economic Review*, 94, 991–1013.
- BOLLINGER, C., AND M. DAVID (2001): “Estimation With Response Error and Nonresponse: Food-Stamp Participation in the SIPP,” *Journal of Business and Economic Statistics*, 19, 129–141.
- BUSSO, M., J. DiNARDO, AND J. MCCRARY (2009): “Finite Sample Properties of Semiparametric Estimators of Average Treatment Effects,” *mimeo*.
- CARROLL, R., D. RUPPERT, AND L. STEFANSKI (1995): *Measurement Error in Nonlinear Models*. Chapman and Hall, London.
- CHEN, P., E. WONG, R. DOMINIK, AND M. STEINER (2000): “A transitional model of barrier methods compliance with unbalanced loss to follow-up,” *Statistics in Medicine*, 19, 71–82.
- CRUMP, R., J. HOTZ, G. IMBENS, AND O. MITNIK (2009): “Dealing with limited overlap in estimation of average treatment effects,” *Biometrika*, 96, 187–199.
- DAS, M., W. NEWWEY, AND F. VELLA (2003): “Nonparametric Estimation of Sample Selection Models,” *Review of Economic Studies*, 70, 33–58.
- DEHEJIA, R., AND S. WAHBA (1999): “Causal Effects in Non-experimental Studies: Reevaluating the Evaluation of Training Programmes,” *Journal of American Statistical Association*, 94, 1053–1062.
- DiNARDO, J., J. MCCRARY, AND L. SANBONMATSU (2006): “Constructive Proposals for Dealing with Attrition: An Empirical Example,” *Working paper, University of Michigan*.
- DING, W., AND S. LEHRER (2010): “Estimating Treatment Effects from Contaminated Multi-period Education Experiments: The Dynamic Impacts of Class Size Reductions,” *Review of Economics and Statistics*, 92, 31–42.
- FITZGERALD, J., P. GOTTSCHALK, AND R. MOFFITT (1998): “An Analysis of Sample Attrition in Panel Data: The Michigan Panel Study of Income Dynamics,” *Journal of Human Resources*, 33, 251–299.
- FRANGAKIS, C., R. BROOKMEYER, R. VARADHAN, M. SAFAEIAN, D. VLAHOV, AND S. STRATHDEE (2004): “Methodology for Evaluating a Partially Controlled Longitudinal Treatment Using Principal Stratification, With Application to a Needle Exchange Program,” *Journal of the American Statistical Association*, 99, 239–249.
- FRANGAKIS, C., AND D. RUBIN (1999): “Addressing complications of intention-to-treat analysis in the combined presence of all-or-none treatment-noncompliance and subsequent missing outcomes,” *Biometrika*, 86, 365–379.
- FRUMENTO, P., F. MEALLI, B. PACINI, AND D. B. RUBIN (2012): “Evaluating the Effect of Training on Wages in the Presence of Noncompliance, Nonemployment, and Missing Outcome Data,” *Journal of the American Statistical Association*, 107, 450–466.

- FRÖLICH, M. (2004): “Finite Sample Properties of Propensity-Score Matching and Weighting Estimators,” *Review of Economics and Statistics*, 86, 77–90.
- (2007): “Nonparametric IV Estimation of Local Average Treatment Effects with Covariates,” *Journal of Econometrics*, 139, 35–75.
- (2008): “Parametric and Nonparametric Regression in the Presence of Endogenous Control Variables,” *International Statistical Review*, 76, 214–227.
- HAUSMAN, J., AND D. WISE (1979): “Attrition Bias In Experimental and Panel Data: The Gary Income Maintenance Experiment,” *Econometrica*, 47(2), 455–473.
- HECKMAN, J. (1976): “The Common Structure of Statistical Models of Truncation, Sample Selection, and Limited Dependent Variables, and a Simple Estimator for such Models,” *Annals of Economic and Social Measurement*, 5, 475–492.
- HECKMAN, J., H. ICHIMURA, AND P. TODD (1997): “Matching as an Econometric Evaluation Estimator: Evidence from Evaluating a Job Training Programme,” *Review of Economic Studies*, 64, 605–654.
- HEITJAN, D., AND S. BASU (1996): “Distinguishing Missing at Random and Missing Completely at Random,” *The American Statistician*, 50, 207–213.
- HERZOG, S. (2005): “Measuring Determinants of Student Return vs. Dropout/Stopout vs. Transfer: A First-to-Second Year Analysis of New Freshmen,” *Research in Higher Education*, 46, 883–928.
- HIRANO, K., G. IMBENS, AND G. RIDDER (2003): “Efficient Estimation of Average Treatment Effects Using the Estimated Propensity Score,” *Econometrica*, 71, 1161–1189.
- HIRANO, K., G. W. IMBENS, G. RIDDER, AND D. B. RUBIN (2001): “Combining Panel Data Sets with Attrition and Refreshment Samples,” *Econometrica*, 69, 1645–1659.
- HORVITZ, D., AND D. THOMPSON (1952): “A Generalization of Sampling without Replacement from a Finite Population,” *Journal of American Statistical Association*, 47, 663–685.
- HUBER, M. (2012): “Identification of average treatment effects in social experiments under alternative forms of attrition,” *Journal of Educational and Behavioral Statistics*, 37, 443–474.
- (2013): “Treatment evaluation in the presence of sample selection,” *forthcoming in Econometric Reviews*.
- HUBER, M., M. LECHNER, AND C. WUNSCH (2013): “The performance of estimators based on the propensity score,” *Journal of Econometrics*, 175, 1–21.
- HUBER, M., AND G. MELLACE (2013): “Testing instrument validity for LATE identification based on inequality moment constraints,” *forthcoming in the Review of Economics and Statistics*.
- IMAI, K. (2008): “Sharp bounds on the causal effects in randomized experiments with ‘truncation-by-death’,” *Statistics & Probability Letters*, 78, 144–149.

- IMBENS, G. (2004): “Nonparametric Estimation of Average Treatment Effects under Exogeneity: A Review,” *Review of Economics and Statistics*, 86, 4–29.
- IMBENS, G., AND J. ANGRIST (1994): “Identification and Estimation of Local Average Treatment Effects,” *Econometrica*, 62, 467–475.
- KITAGAWA, T. (2013): “A Bootstrap Test for Instrument Validity in the Heterogeneous Treatment Effect Model,” *UCL working paper*.
- KYRIAZIDOU, E. (1997): “Estimation of a Panel Data Sample Selection Model,” *Econometrica*, 65, 1335–1364.
- (2001): “Estimation of Dynamic Panel Data Sample Selection Models,” *The Review of Economic Studies*, 68, 543–572.
- LECHNER, M. (1999): “Earnings and Employment Effects of Continuous Off-the-Job Training in East Germany after Unification,” *Journal of Business and Economic Statistics*, 17, 74–90.
- LECHNER, M., R. MIQUEL, AND C. WUNSCH (2011): “Long-Run Effects of Public Sector Sponsored Training in West Germany,” *Journal of the European Economic Association*, 9, 742–784.
- LEE, D. (2009): “Training, Wages, and Sample Selection: Estimating Sharp Bounds on Treatment Effects,” *Review of Economic Studies*, 76, 1071–1102.
- LEPPEL, K. (2002): “Similarities and differences in the college persistence of men and women,” *Review of Higher Education*, 25, 433–450.
- LITTLE, R., AND D. RUBIN (1987): *Statistical Analysis with Missing Data*. Wiley, New York.
- LITTLE, R. J. A. (1995): “Modeling the Drop-Out Mechanism in Repeated-Measures Studies,” *Journal of the American Statistical Association*, 90, 1112–1121.
- LOK, J., R. GILL, A. VAN DER VAART, AND J. ROBINS (2004): “Estimating the causal effect of a time-varying treatment on time-to-event using structural nested failure time models,” *Statistica Neerlandica*, 58, 271–295.
- MACKINNON, J. G. (2006): “Bootstrap Methods in Econometrics,” *The Economic Record*, 82, S2–S18.
- MANSKI, C. (1989): “Anatomy of the Selection Problem,” *Journal of Human Resources*, 24, 343–360.
- (1990): “Nonparametric Bounds on Treatment Effects,” *American Economic Review, Papers and Proceedings*, 80, 319–323.
- MATTEI, A., AND F. MEALLI (2007): “Application of the Principal Stratification Approach to the Faenza Randomized Experiment on Breast Self-Examination,” *Biometrics*, 63, 437–446.
- MEALLI, F., G. IMBENS, S. FERRO, AND A. BIGGERI (2004): “Analyzing a randomized trial on breast self-examination with noncompliance and missing outcomes,” *Biostatistics*, 5, 207–222.
- MURPHY, S., M. VAN DER LAAN, AND J. ROBINS (2001): “Marginal mean models for dynamic regimes,” *Journal of the American Statistical Association*, 96, 1410–1423.

- NEWWEY, W. (1994): “The Asymptotic Variance of Semiparametric Estimators,” *Econometrica*, 62, 1349–1382.
- PENG, Y., R. J. A. LITTLE, AND T. E. RAGHUNATHAN (2004): “An Extended General Location Model for Causal Inferences from Data Subject to Noncompliance and Missing Values,” *Biometrics*, 60, 598–607.
- POWELL, J., J. STOCK, AND T. STOKER (1989): “Semiparametric Estimation of Index Coefficients,” *Econometrica*, 57, 1403–1430.
- PREISSER, J. S., A. T. GALECKI, K. K. LOHMAN, AND L. E. WAGENKNECHT (2000): “Analysis of Smoking Trends with Incomplete Longitudinal Binary Responses,” *Journal of the American Statistical Association*, 95, 1021–1031.
- RACINE, J., AND Q. LI (2004): “Nonparametric Estimation of Regression Functions with Both Categorical and Continuous Data,” *Journal of Econometrics*, 119, 99–130.
- ROBINS, J. (1989): “The Analysis of Randomized and Non-Randomized AIDS Treatment Trials Using a New Approach to Causal Inference in Longitudinal Studies,” in *Health Service Research Methodology: A Focus on AIDS*, ed. by L. Sechrest, H. Freeman, and A. Mulley, pp. 113–159. U.S. Public Health Service, Washington, DC.
- ROBINS, J., S. GREENLAND, AND F.-C. HU (1999): “Estimation of the causal effect of a time-varying exposure on the marginal mean of a repeated binary outcome,” *Journal of the American Statistical Association*, 94, 687–700.
- ROBINS, J., A. ROTNITZKY, AND L. ZHAO (1995): “Analysis of Semiparametric Regression Models for Repeated Outcomes in the Presence of Missing Data,” *Journal of American Statistical Association*, 90, 106–121.
- ROBINS, J. M. (1986): “A new approach to causal inference in mortality studies with a sustained exposure period application to control of the healthy worker survivor effect,” *Mathematical Modelling*, 7, 1393–1512.
- ROBINS, J. M., A. ROTNITZKY, AND L. ZHAO (1994): “Estimation of Regression Coefficients When Some Regressors Are not Always Observed,” *Journal of the American Statistical Association*, 90, 846–866.
- ROSENBAUM, P., AND D. RUBIN (1983): “The Central Role of the Propensity Score in Observational Studies for Causal Effects,” *Biometrika*, 70, 41–55.
- ROTNITZKY, A., J. ROBINS, AND D. SCHARFSTEIN (1998): “Semiparametric Regression for Repeated Outcomes With Nonignorable Nonresponse,” *Journal of the American Statistical Association*, 93, 1321–1339.
- RUBIN, D. (1976): “Inference and Missing Data,” *Biometrika*, 63, 581–592.
- (1977): “Formalizing Subjective Notions About the Effect of Nonrespondents in Sample Surveys,” *Journal of the American Statistical Association*, 72, 538–543.

- (1978): “Multiple Imputations in Sample Surveys-A Phenomenological Bayesian Approach to Nonresponse,” in *Proceedings of the Survey Research Methods Section, American Statistical Association*, pp. 20–34.
- SCHARFSTEIN, D., A. ROTNITZKY, AND J. ROBINS (1999): “Adjusting for Nonignorable Drop-Out Using Semiparametric Nonresponse Models,” *Journal of the American Statistical Association*, 94, 1096–1120.
- SEMYKINA, A., AND J. WOOLDRIDGE (2006): “Estimating Panel Data Models in the Presence of Endogeneity and Selection: Theory and Application,” *unpublished manuscript*.
- SERFLING, R. (1980): *Approximation Theorems of Mathematical Statistics*. Wiley, New York.
- SHAH, A., N. LAIRD, AND D. SCHOENFELD (1997): “A Random-Effects Model for Multiple Characteristics With Possibly Missing Data,” *Journal of the American Statistical Association*, 92, 775–779.
- SHEPHERD, B. E., M. W. REDMAN, AND D. P. ANKERST (2008): “Does Finasteride Affect the Severity of Prostate Cancer? A Causal Sensitivity Analysis,” *Journal of the American Statistical Association*, 103, 1392–1404.
- TAN, Z. (2006): “Regression and Weighting Methods for Causal Inference Using Instrumental Variables,” *Journal of the American Statistical Association*, 101, 1607–1618.
- TINTO, V. (1997): “Classroom as communities: Exploring the educational character of student persistence,” *Journal of Higher Education*, 68, 599–623.
- VAN DER LAAN, M., AND D. RUBIN (2006): “Targeted Maximum Likelihood Learning,” *The International Journal of Biostatistics*, 2.
- WANG, L., A. ROTNITZKY, X. LIN, R. E. MILLIKAN, AND P. F. THALL (2012): “Evaluation of Viable Dynamic Treatment Regimes in a Sequentially Randomized Trial of Advanced Prostate Cancer,” *Journal of the American Statistical Association*, 107, 493–508.
- XIE, H., AND Y. QIAN (2012): “Measuring the impact of nonignorability in panel data with non-monotone nonresponse,” *Journal of Applied Econometrics*, 27, 129–159.
- YAU, L., AND R. LITTLE (2001): “Inference for the complier-average causal effect from longitudinal data subject to noncompliance and missing data, with application to a job training assessment for the unemployed,” *Journal of American Statistical Association*, 96, 1232–1244.
- ZHANG, J., AND D. RUBIN (2003): “Estimation of causal effects via principal stratification when some outcome are truncated by death,” *Journal of Educational and Behavioral Statistics*, 28, 353–368.
- ZHANG, J., D. RUBIN, AND F. MEALLI (2008): “Evaluating The Effects of Job Training Programs on Wages through Principal Stratification,” in *Advances in Econometrics: Modelling and Evaluating Treatment Effects in Econometrics*, ed. by D. Millimet, J. Smith, and E. Vytlacil, vol. 21, pp. 117–145. Elsevier Science Ltd.
- (2009): “Likelihood-Based Analysis of Causal Effects of Job-Training Programs Using Principal Stratification,” *Journal of the American Statistical Association*, 104, 166–176.