

Gautier, Eric; Hoderlein, Stefan

**Working Paper**

## A triangular treatment effect model with random coefficients in the selection equation

cemmap working paper, No. CWP39/12

**Provided in Cooperation with:**

The Institute for Fiscal Studies (IFS), London

*Suggested Citation:* Gautier, Eric; Hoderlein, Stefan (2012) : A triangular treatment effect model with random coefficients in the selection equation, cemmap working paper, No. CWP39/12, Centre for Microdata Methods and Practice (cemmap), London, <https://doi.org/10.1920/wp.cem.2012.3912>

This Version is available at:

<https://hdl.handle.net/10419/79537>

**Standard-Nutzungsbedingungen:**

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

**Terms of use:**

*Documents in EconStor may be saved and copied for your personal and scholarly purposes.*

*You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.*

*If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.*

# A triangular treatment effect model with random coefficients in the selection equation

---

**Eric Gautier**  
**Stefan Hoderlein**

The Institute for Fiscal Studies  
Department of Economics, UCL

**cemmap** working paper CWP39/12

# A TRIANGULAR TREATMENT EFFECT MODEL WITH RANDOM COEFFICIENTS IN THE SELECTION EQUATION

ERIC GAUTIER AND STEFAN HODERLEIN

ABSTRACT. In this paper we study nonparametric estimation in a binary treatment model where the outcome equation is of unrestricted form, and the selection equation contains multiple unobservables that enter through a nonparametric random coefficients specification. This specification is flexible because it allows for complex unobserved heterogeneity of economic agents and non-monotone selection into treatment. We obtain conditions under which both the conditional distributions of  $Y_0$  and  $Y_1$ , the outcome for the untreated, respectively treated, given first stage unobserved random coefficients, are identified. We can thus identify an average treatment effect, conditional on first stage unobservables called UCATE, which yields most treatment effects parameters that depend on averages, like ATE and TT. We provide sharp bounds on the variance, the joint distribution of  $(Y_0, Y_1)$  and the distribution of treatment effects. In the particular case where the outcomes are continuously distributed, we provide novel and weak conditions that allow to point identify the joint conditional distribution of  $Y_0, Y_1$ , given the unobservables. This allows to derive every treatment effect parameter, *e.g.* the distribution of treatment effects and the proportion of individuals who benefit from treatment. We present estimators for the marginals, average and distribution of treatment effects, both conditional on unobservables and unconditional, as well as total population effects. The estimators use all the data and discard tail values of the instruments when they are too unlikely. We provide their rates of convergence, and analyze their finite sample behavior in a simulation study. Finally, we also discuss the situation where some of the instruments are discrete.

---

*Date:* First Version: September 2010 ; This version: November 29, 2012.

*Keywords:* Treatment Effects, Endogeneity, Random Coefficients, Nonparametric Identification, Partial Identification, Roy Model, Ill-Posed Inverse Problems, Deconvolution, Radon Transform, Rates of Convergence.

We are grateful to seminar participants at Boston College, Chicago, CREST, Harvard-MIT, Kyoto, Nanterre, Northwestern, Oxford, Princeton, Toulouse, UCL, Vanderbilt, 2011 CIRM New Trends in Mathematical Statistics, 2012 Bates White, CLAPEM, ESEM, SCSE, and the World Congress in Probability and Statistics conferences and Arnaud Debussche for helpful comments. The authors are very grateful to Helen Broome, who is coauthor of a companion paper considering the evaluation of the returns to college education, for her many useful remarks and assistance in the simulation study.

## 1. INTRODUCTION

In this paper we consider estimation of treatment effect parameters in the presence of multiple sources of unobserved heterogeneity in the selection equation. We consider the following treatment effect model

$$(1.1) \quad Y = Y_0 + \Delta D, \quad \text{where } \Delta = Y_1 - Y_0,$$

$$(1.2) \quad D = \mathbf{1} \{V - Z'\Gamma - \Theta > 0\}.$$

The outcome equation (1.1) is a linear model with a binary endogenous regressor  $D$  and random coefficients where,  $Y_0$  is the random intercept and  $\Delta$  the random slope. Note that this outcome equation allows for unrestricted heterogeneity, since it is equivalent to a nonseparable model  $Y = \psi(D, U)$ , with  $U$  being a (potentially infinite) vector of unobserved variables. The random coefficients have an interpretation:  $Y_0$  is the outcome in the control group or base state,  $Y_1$  the outcome in the treated group, and  $\Delta = Y_1 - Y_0$  is the net gain from an ideal exogenous move of an individual from state 0 to state 1, called the effect of treatment. In our model,  $Y_0$  and  $Y_1$  may be continuous or discrete, and individuals are observed in either of the two states 0 or 1. We aim to estimate features of the random slope  $\Delta$  (for example its average is the average treatment effect (ATE)) or the joint distribution of  $(Y_0, \Delta)$  which yields the joint distribution of potential outcomes  $(Y_0, Y_1)$ . In this model, the binary regressor  $D$  is endogenous because participation in the treatment is endogenous. We therefore supplement model (1.1) by modeling explicitly the regressor  $D$  *i.e.*, the selection into treatment. Individuals select themselves into treatment, if the net (expected) utility of participating in the treatment is positive, as formalized in equation (1.2), where  $\mathbf{1}$  denotes the indicator function. This net utility depends on a vector of instrumental variables  $(V, Z') \in \mathbb{R}^L$  which are observable to the econometrician. It also depends, in a nonseparable fashion, on unobserved parameters  $(\Gamma', \Theta) \in \mathbb{R}^L$  which are allowed to vary across the population. Because the scale of the net utility cannot be identified, we set the coefficient of  $V$  to 1. This can be done if, in the original scale, the coefficient of  $V$  is positive<sup>1</sup>.

To fix ideas, one may think of the instruments as cost factors or elements of information about the net utility of treatment, and of the random slope coefficients as reflecting the heterogeneous impact of these factors on net utility. The random intercept can be interpreted as contributions to net utility that are unobserved to the econometrician such as the anticipated gains of being treated plus

---

<sup>1</sup>We can change  $V$  in  $-V$  if it is negative.

possibly some random term (*e.g.*, the random intercept of the cost function). While we allow for a rich structure in terms of unobservables that goes significantly beyond the common “scalar unobservable threshold crossing” model, we remark that at the same time we place potentially restrictive structure by requiring that this model be linear in parameters, and that the unobservables in the selection equation have the same dimension as the vector of exogenous variables. While we do not literally believe in an exact linear structure on individual level, we think of a linear index structure as a good first order approximation. Rather than aiming to capture higher order terms in observable variables, in this paper we want to place emphasis on the dependence of the participation decision on an unobserved structure in a way that is more in line with structural economics.

The main results in this paper establish that under very general conditions,

$$(1.3) \quad f_{\Gamma', \Theta}, f_{Y_j | \Gamma', \Theta}^2, j = 0, 1, \quad \text{and} \quad \mathbb{E}[\Delta | \Gamma = \gamma, \Theta = \theta],$$

are point identified. Moreover, we provide sharp bounds on  $Var[\Delta]$  and  $F_\Delta$ . Finally when the outcome is continuous, under additional conditions we show that  $Var[\Delta | \Gamma = \gamma, \Theta = \theta]$ ,  $f_{\Delta | \Gamma, \Theta}$  and  $f_{Y_0, Y_1 | \Gamma, \Theta}$  are point identified. Let us now elaborate on the individual objects.

The average treatment effect, conditional on first stage preference parameters,  $\mathbb{E}[\Delta | \Gamma = \gamma, \Theta = \theta]$ , abbreviated UCATE, is similar in spirit to the marginal treatment effect (MTE, see Björklund and Moffitt (1987), and Heckman and Vytlacil (2005)) and shares many of its appealing properties (policy invariance, interpretation in terms of willingness to pay for people at the margin of indifference, averages like the ATE are straightforwardly derived, etc.).

It differs in as far as instead of depending on a single first stage unobservable, it allows to condition on the entire vector  $\Gamma, \Theta$  of heterogenous first stage parameters. Unlike the scalar unobservable threshold crossing model, the selection equation neither imply monotonicity nor uniformity (see Imbens and Angrist (1994) and Heckman and Vytlacil (2005)), as soon as  $L \geq 2$ . It thus allows for more general heterogeneity patterns in the selection equation. In particular, there may be both compliers and defiers in the population. Because of this generality, the model (1.2) is suggested in Heckman and Vytlacil (2005) as a benchmark nonseparable, nonmonotonic model.

The marginals,  $f_{Y_j | \Gamma', \Theta} f_{\Gamma', \Theta}$ ,  $j = 0, 1$ , are identified under the same conditions as UCATE without appealing to randomized experiments or selection on unobservables. From the marginals we can recover the unconditional marginals, as well as many other parameters, like the quantile treatment

---

<sup>2</sup>Throughout this paper, we will refer to the density and cumulative distribution function (CDF for short) of a vector  $A$  as  $f_A$ , respectively  $F_A$ , we will write  $f_{A|B}(\cdot|b)$  and  $F_{A|B}(\cdot|b)$  the conditional densities and CDF of  $A$  given  $B = b$ .

effect (QTE), see Abadie, Angrist and Imbens (2000) and Chernozhukov and Hansen (2005). We also derive sharp bounds for (1) the variance of treatment effect (VATE), (2) the CDF of the two potential outcomes, and (3) the CDF of treatment effects.<sup>3</sup>

These bounds are potentially wide, and obtaining point identification under plausible conditions is thus desirable. To point identify the conditional variance of treatment effects (UCVATE)  $Var[\Delta|\Gamma = \gamma, \Theta = \theta]$ , and thus also the unconditional variance, we impose the additional assumption that conditional on  $(\Gamma', \Theta)$ ,  $Y_0$  and  $\Delta$  are uncorrelated. For the Unobservables Conditioned Distribution of Treatment Effects (UCDITE)  $f_{\Delta|\Gamma', \Theta}$  we impose the stronger assumption that  $Y_0$  and  $\Delta$  are independent conditional on  $\Gamma, \Theta$  (which we denote as  $Y_0 \perp \Delta|\Gamma, \Theta$ )<sup>4</sup>. Conditional independence assumptions between the gain  $\Delta$  and the base state  $Y_0$  are made in Heckman, Smith and Clements (1997), Heckman and Clements (1998). However, in these references the independence assumption is conditional on  $D$ , or on observable variables  $X$ . In contrast, one attractive feature of the approach put forward in this paper is that the independence is conditional on the unobservables entering the selection equation (as well as control variables  $X$ ). This makes this assumption more likely to hold, as we argue in detail in Section 3.4 using extensions of the Roy model. The unobservables in the selection equation contain information on ex-ante forecast of the gains and cost factors. At this point, we would only like to point out that this assumption is satisfied, if there exist otherwise unrestricted mappings  $\psi_0, \psi_\Delta$  such that  $Y_0 = \psi_0(\Gamma', \Theta, U_0)$ , and  $\Delta = \psi_\Delta(\Gamma', \Theta, U_\Delta)$ , with  $U_0, U_\Delta$  possibly infinite dimensional, such that  $U_0 \perp U_\Delta \perp (\Gamma', \Theta)$ .<sup>5</sup> This paper therefore shows that allowing for several sources of unobserved heterogeneity in the selection equation is important, not just in its own right,

---

<sup>3</sup>The bounds for (2) stem from the classical bounds of Hoeffding (1940) and Frechet (1951) and are used in Heckman, Smith and Clements (1997), Manski (1997) and Heckman and Smith (1998). Fan and Park (2010) and Firpo and Ridder (2008) apply the Makarov (1981) bounds to infer bounds on the distribution of treatment effects, Fan and Zhu (2009) obtain bounds for functionals of the joint distribution of potential outcomes, like inequality measures based on the distribution of treatment effects. Unlike these references the bounds are obtained in the case of possible selection on observables and unobservables. As we shall see, the bounds on the distribution of treatment effects are also sharper than what a direct application of the Makarov bounds would yield because they are averages of Makarov bounds on CDFs of the marginals conditional on the observables and unobservables of equation (1.2).

<sup>4</sup>We indeed assume  $Y_0 \perp \Delta|\Gamma, \Theta, X$  with some observables  $X$  that are used as control variables. This can make the independence assumption more likely in the same spirit as the missing at random assumption in the missing data literature or simply allow inferring treatment effects for population subgroups.

<sup>5</sup>This can be interpreted as the fact that the selection equation reveals information about the common endogenous factors; there is (potentially complicated) remaining heterogeneity in  $Y_0$  and  $\Delta$ , but it is independent of everything else.

but also to allow for more dependence between the gain  $\Delta$  and the base state  $Y_0$ , and still point identify distributional treatment effect parameters.

Another material assumption is a support condition. When  $L$ , the dimension of instruments is 2 or larger, we impose a condition relating the support of the instruments and that of the unobserved heterogeneity parameters. Though we can deal with instruments with bounded support, it is a type of “large” support assumption which is required to fully recover the entire distribution of random coefficients in the selection equation. When  $L = 1$ , we establish that when the variation of our instruments is small relative to that of the unobserved heterogeneity parameters, we can identify the average, variance or distribution of effects for the subpopulation defined by the range of the instrument. This subpopulation is related to the one considered in Angrist, Graddy and Imbens (2000), the LATE of Imbens and Angrist (1994) being a special case. We suspect that something similar is feasible for  $L \geq 2$ , but leave it for future research.

Based on these identification results, we provide sample counterparts estimators and obtain their rates of convergence. It is known from Beran, Feuerwerker and Hall (1996), Hoderlein, Klemelä and Mammen (2010) and Gautier and Kitamura (2009) that the estimation of the distribution of random coefficients in the single equation case with exogenous regressors is an ill-posed inverse problem. We extend these papers by allowing dependence between the random coefficients and the regressors and by relaxing the full support assumption on the regressors<sup>6</sup>. Similar to Imbens and Newey (2009), we deal with an endogenous regressor  $D$  in equation (1.1) through a triangular system with an equation for the endogenous regressor but we allow for multiple sources of unobserved heterogeneity entering in a nonseparable fashion in (1.2).

In contrast to all of these references, in our approach the first stage (selection) equation allows for multiple sources of unobserved heterogeneity, and monotonicity is not imposed. Estimation of the marginals  $(Y_0 + \Delta, \Gamma', \Theta)$ ,  $(Y_0, \Gamma', \Theta)$  and  $(\Gamma', \Theta)$ , UCATE, VATE or extensions of the QTE, then relies on solving ill-posed inverse problems that involve the Radon transform. The estimation of the distribution of treatment effects is a deconvolution problem and features thus in addition another ill-posed inverse problem (see also Heckman, Smith and Clements (1997) and Heckman and Smith (1998)). More precisely, in our setup, it consists of a conditional deconvolution problem with unknown but estimable distributions of (1) the signal plus error and of (2) the error. Evdokimov (2010) considers conditional deconvolution in nonparametric panel data models with unobserved heterogeneity. In

---

<sup>6</sup>In (1.1)  $D$  and  $(Y_0, \Delta)$  are dependent and  $D$  has a limited support, in (1.2) we will allow  $V$  to have limited support and  $(V, Z')$  and  $(\Gamma', \Theta)$  to be dependent but independent given control variables.

classical deconvolution, the density of  $Y_0$  is known and the characteristic function of  $Y_0 + \Delta$  is estimated via the empirical characteristic function which estimates the true characteristic function at rate  $1/\sqrt{N}$ . An extension studied in the statistics literature considers the case where the density of  $Y_0$  is estimable at rate  $1/\sqrt{N}$  using a preliminary sample (see, *e.g.* Neumann (1997), Johannes (2009) and Comte and Lacour (2011)). In this paper the Fourier transforms of the densities of  $(Y_0 + \Delta, \Gamma', \Theta)$  and  $(Y_0, \Gamma', \Theta)$  with respect to the first argument are estimated solving inverse problems using the same sample. Getting unconditional parameters requires to compute an integral over the conditional effects weighted by the joint density of the random coefficients.

More generally, this paper touches upon two related sets of literatures. The first is the treatment effect literature, in particular the part that is related to distributional treatment effects, the second is the random coefficients literature. First, we obtain results for treatment effect parameters that depend on averages. This is related to the important contributions of LATE (Imbens and Angrist (1994)) and MTE (Björklund and Moffitt (1987), Heckman and Vytlacil (1999, 2005, 2007)). But we also obtain results on the marginals which are related to the quantile treatment effects of Abadie, Angrist and Imbens (2002), Chernozhukov and Hansen (2005), and to results by Heckman, Smith and Clements (1997), Heckman and Smith (1998), Carneiro, Hansen and Heckman (2003), Aakvik, Heckman and Vytlacil (2005), Abbring and Heckman (2007) and Fan and Zhu (2009), Fan and Park (2010) and Firpo and Ridder (2008). Note that the first two results on quantile treatment effect essentially require a rank invariance assumption, *i.e.*, the individuals retain their ordering both in the treatment and the control group, an assumption which may only be weakened slightly. This assumption is restrictive, and has been criticized, see Heckman, Smith and Clements (1997). Carneiro, Hansen and Heckman (2003) and Aakvik, Heckman and Vytlacil (2005) consider a factor model approach that allows to identify the distribution of treatment effects. Our identification strategy has some similarities with that of Lewbel (2000), which considers the estimation of the average of the random coefficients in a binary choice equation, and Lewbel (2007), which considers a selection model. However, here we recover both the distribution of the vector of random coefficient in the binary choice and the whole distribution of potential outcomes. Moffitt (2008) considers a model with multiple sources of unobserved heterogeneity in the selection equation, which enter in a more general form than the linear index structure of equation (1.2), but he imposes monotonicity. Klein (2010) considers a different specification where a second independent source of unobserved heterogeneity enters the scalar unobservable threshold crossing in a nonseparable fashion.



The second related line of work is the literature on nonparametric random coefficients models. Random coefficients models allow the preference or production parameters to vary across the population. In this paper, we specifically allow for different individuals to have different costs and benefits of treatment. We emphasize the nonparametric aspect of our analysis, which allows to be flexible about the form of unobserved heterogeneity. References in econometrics include Elbers and Ridder (1982), Heckman and Singer (1984), Beran and Hall (1992), Beran, Feuerverger and Hall (1996), Ichimura and Thompson (1998), Fox and Gandhi (2011), Hoderlein, Klemelä and Mammen (2010), Gautier and Kitamura (2009). Gautier and Le Pennec (2011) obtain minimax lower bounds and an adaptive data-driven estimator for the estimation of the distribution of random coefficients in random coefficients binary choice models. The last three references focus on estimation and continuous random coefficients, and recognize that the estimation of the density of the latent random coefficients vector is a statistical inverse problem in this scenario. The literature on the treatment of these problems is extensive in statistics and econometrics, and we refer to Carrasco, Florens and Renault (2007) for a survey of applications in economics. Fox and Gandhi (2009) study identification of the distribution of unobserved heterogeneity in Roy models, however they focus on the case of discretely supported random coefficients, and do not allow for a random intercept in the selection equation.

This paper is structured as followed. In section 2, we introduce the model formally, and discuss the basic assumptions. In section 3 we establish the main identification results. We show that under the baseline assumptions the marginals of each potential outcome are identified, conditional on random coefficients. This allows to obtain the UCATE, many other treatment effect parameters that only depend on averages, and other parameters that solely depend on marginals, like the QTE. Moreover, we obtain bounds on the variance of treatment effects, on the joint distribution of the two potential outcomes, and consequently also on the distribution of treatment effects. Finally, in section 3.4, we introduce two nested assumptions that allow to point identify the variance and the distribution of treatment effects. We discuss the interpretation of these assumptions and their relevance in the context of extensions of the Roy model, discuss how they aid in identification, and present again sample counterparts estimators. In section 4, we present estimators and obtain their rates of convergence. The estimators are not based on theoretical formulas that only involve values at “infinity” of the instruments in the selection equation, nor situations with unselected samples. It is rather the contrary, as they make efficient use of all the data and involve trimming of tail values of the instruments. In section 5, we analyze the finite sample behavior through a simulation study. In section 6 we consider an

alternative scaling of the random coefficients binary choice and the case where some of the instruments are binary.

## 2. THE RANDOM COEFFICIENTS MODEL AND ASSUMPTIONS

This section introduces the formal setup in which we analyze the effects of treatment. We distinguish between two cases. The first one is the case where we have a vector of unobservables, and is the core innovation in this paper. The second one is the “traditional” case, which features a scalar unobservable and is displayed largely for comparison. Before we discuss these scenarios in detail, we start, however, by introducing some crucial pieces of notation and basic probabilistic assumptions.

Throughout this paper, we use uppercase letters for random variables and lowercase for their realizations. In addition to the observable variables  $(Y, D, V, Z)$  which have already been introduced above, we assume that there might be another observable random vector  $X$  on which we may want to condition upon when doing inference. Examples include household characteristics like age, gender, race etc. We do not impose any restriction on the dependence of  $Y_0, Y_1, V, Z, \Gamma, \Theta$  on  $X$ , besides regularity conditions which we detail below. As with  $Y$ , we denote by  $X_0$  and  $X_1$  the random variable  $X$  when  $D$  is 0, respectively 1.

The data consists of the realizations of  $N$  independent and identically distributed copies of the population random variables, which we denote as  $(y_i, d_i, v_i, z'_i, x'_i)_{i=1, \dots, N}$ , where  $N$  is the sample size. We denote by  $\text{supp}(A)$  or  $\text{supp}(f_A)$ , the support of the random vector  $A$  and by  $\text{Int}(A)$  the interior of a set  $A$ .

For mathematical convenience, in the case where  $L \geq 2$ , it is useful to renormalize the index in (1.2). This could be done in several ways. In this paper, we assume that the (random) coefficient of  $V$  in the original net utility scale has a sign. A more general sufficient condition is presented in Section 6.1. The approach put forward in this paper allows us to handle the case where  $V$  has bounded support, too. As a next step, we divide by  $\|(Z', 1)\|$ , and use the notations  $\tilde{S} = (Z', 1)' / \|(Z', 1)\|$ ,  $\tilde{V} = V / \|(Z', 1)\|$  and  $\tilde{\Gamma} = (\Gamma', \Theta)'$ . Then, (1.2) becomes

$$(2.1) \quad D = \mathbf{1}\{\tilde{S}'\tilde{\Gamma} < \tilde{V}\}.$$

We invoke the following assumptions (when  $L = 1$  simply drop  $Z$  and  $\Gamma$  below).

**Assumption 2.1.** (A-1) The conditional distribution of  $(\tilde{S}', \tilde{V}, \Theta, \Gamma')$  given  $X = x$  is absolutely continuous with respect to the product of the spherical measure on  $\mathbb{S}^{L-1}$  and the Lebesgue measure on  $\mathbb{R}^{L+1}$  for almost every  $x$  in  $\text{supp}(X)$  ;

(A-2)  $(V, Z) \perp (Y_0, \Gamma', \Theta) | X$  and  $(V, Z) \perp (Y_1, \Gamma', \Theta) | X$  ;

(A-3)  $0 < \mathbb{P}(D = 1 | X) < 1$  *a.s.* ;

(A-4)  $X_0 = X_1$  *a.s.*

Assumption (A-1) defines the setup of this paper as one with continuous instruments. Note that the fact that  $(\tilde{S}', \tilde{V})$  is continuous does not require that in the original scale  $(Z', V)$  be continuous. For example, it is possible that  $V$  is binary and  $\tilde{V}$  is continuous. We will see in Section 6.2 how to handle other binary instruments. In Sections 3.4.3, 4.3.3 and 4.3.4, we strengthen (A-1) to

(A-1') The conditional distribution of  $(\tilde{S}', \tilde{V}, Y_0, Y_1, \Theta, \Gamma')$  given  $X = x$  is absolutely continuous with respect to the product of the spherical measure on  $\mathbb{S}^{L-1}$  and the Lebesgue measure on  $\mathbb{R}^{L+3}$  for almost every  $x$  in  $\text{supp}(X)$ .

It is only in these sections that we consider the particular case where the outcomes are continuous. The responses  $Y_0$  and  $Y_1$  are allowed to be heterogeneous in a general way and of the form  $Y_0 = \psi_0(X, \Gamma', \Theta, U_0)$  and  $Y_1 = \psi_1(X, \Gamma', \Theta, U_1)$  where  $U_0$  and  $U_1$  account for unobserved heterogeneity and can be infinite dimensional. Assumption (A-1) also implies an exclusion restriction: conditional on  $X = x$ ,  $(\tilde{S}, \tilde{V})$  is continuous and has variation. In practice, this is achieved when  $(V, Z')$  are not in the list of regressors  $X$ .

Assumption (A-2) requires the instruments to be independent of the random parameters given some variables  $X$  which are either exogenous or act as control variables. We allow  $(Y_0, \Gamma', \Theta)$  and  $(Y_1, \Gamma', \Theta)$  to depend on  $(V, Z)$  (unconditional endogeneity of the instruments) as long as we have at hand control variables  $X$  which yield independence. Randomization or pseudorandomization is a classical tool to generate instrumental variables satisfying this assumption unconditional on  $X$ . When  $L \geq 2$ , Assumption (A-2) can be written in terms of the renormalized instruments as

$$(\tilde{V}, \tilde{S}) \perp (Y_0, \tilde{\Gamma}') | X \quad \text{and} \quad (\tilde{V}, \tilde{S}) \perp (Y_1, \tilde{\Gamma}') | X.$$

Assumption (A-3) states that for any  $x \in X$ , there is always a fraction of the population that participates in treatment, and one that does not. Finally, Assumption (A-4) states that  $X$  is not caused by the treatment. As in Heckman and Vytlacil (2005), this last assumption is not strictly required for econometric analysis. However, it makes the inferred quantities more interpretable, and allows to still capture the total effects of  $D$  on  $Y$ , after conditioning on  $X$ .

It is possible, when considering only one unobservable, to consider a model where (1.2) is replaced by

$$(2.2) \quad D = \mathbf{1}\{\mu(V, Z) > \Theta\}$$

where  $\mu$  is a CDF and  $\Theta|X \sim \mathcal{U}(0, 1)$ . This is a well established model studied, among others, in Heckman and Vytlacil (1998, 2005, 2007). We do not present the extension of their results in terms of variance and distribution of treatment effects in this text in order to make the presentation more concise. The arguments are closely related to the single unobservable case used for purpose of comparison below. Already with one unobservable, our results on the variance and distribution of treatment effects are new. The one unobservable in the selection equation framework extends the IV framework where general heterogeneity in response but not in choices is allowed. Additive separability in a scalar unobservable has a strong implication called “monotonicity” in Imbens and Angrist (1994) and “uniformity” in Heckman and Vytlacil (2005): conditional on  $X = x$  where  $x$  belongs to  $\text{supp}(X)$ , for any  $(v, z)$  and  $(v', z')$  in  $\text{supp}(V, Z')$ , if the instruments are moved for everyone from  $(v, z)$  to  $(v', z')$  then either for every  $\theta \in [0, 1]$ ,  $\mathbf{1}\{\mu(v, z) > \theta\} \geq \mathbf{1}\{\mu(v', z') > \theta\}$  or for every  $\theta \in [0, 1]$ ,  $\mathbf{1}\{\mu(v, z) > \theta\} < \mathbf{1}\{\mu(v', z') > \theta\}$ , *i.e.*, there are either no compliers, or no defiers. This is a substantial restriction regarding heterogeneity in choices of treatment (see, *e.g.*, Imbens and Angrist (1994) for several examples). Vytlacil 2002 shows that model (2.2) is equivalent to the monotonicity assumption. In Heckman and Vytlacil (2005) total population average effects are obtained as weighted integrals of the MTE over the whole range of values of the score from 0 to 1. The bounds for integration 0 and 1 correspond to situations without sample selection. This feature is due to monotonicity. The case with more than two sources of unobserved heterogeneity in  $D$  that we advocate in this paper allows for more general heterogeneity in treatment choices in which monotonicity breaks down. The model that we consider is also not additively separable in the components of unobserved heterogeneity in  $\Gamma$ . Heckman and Vytlacil (2005) qualifies the model studied in this paper as the benchmark nonseparable, nonmonotonic model. Introducing multiple sources of unobserved heterogeneity in the selection equation to relax monotonicity is well motivated, see, *e.g.* Heckman and Vytlacil (2005) and Klein (2010). To see why model (1.2) does not impose monotonicity, fix  $v \in \text{supp}(\tilde{V})$  and take  $s$  and  $s'$  in  $\text{supp}(f_{\tilde{S}|\tilde{V}}(\cdot|v))$  and denote by  $D_s(\gamma) = \mathbf{1}\{s'\gamma < v\}$ . In Figure 1 we consider the case where  $L = 2$ ,  $v = 0$  and thus the origin is where the two lines (defined through their normal vectors  $s$  and  $s'$ ) intersect. For an unobserved heterogeneity parameter  $\gamma$  in zone 2,  $D_s = 0$  and  $D_{s'} = 1$  while, for  $\gamma$  in zone 4,  $D_s = 1$  and  $D_{s'} = 0$ , *i.e.*, parts of the population may

be compliers, parts defiers. It is obvious that model (1.2) does not imply unselected samples in the limit if one of the component of  $Z$  goes to infinity. This is related to the fact that monotonicity is not any longer required to hold.

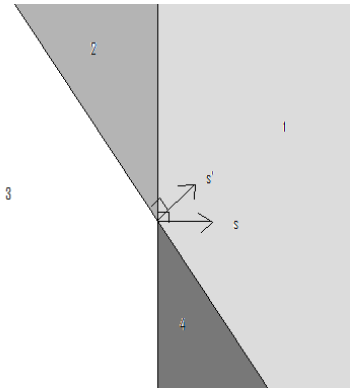


FIGURE 1. Non uniformity with more than 2 unobservables

Assumptions (A-1) and (A-2) together yield that the score function satisfies

$$(2.3) \quad \pi(v, x) = \int_{-\infty}^v f_{\Theta|X}(t|x) dt$$

when  $L = 1$ . In the case when  $L \geq 2$ ,

$$(2.4) \quad \pi(s, v, x) = \mathbb{P}(D = 1 | \tilde{S} = s, \tilde{V} = v, X = x)$$

$$(2.5) \quad = \int_{\mathbb{R}^L} \mathbf{1}\{\gamma' s < v\} f_{\tilde{\Gamma}|X}(\gamma|x) d\gamma$$

$$(2.6) \quad = \int_{-\infty}^v \int_{P_{s,u}} f_{\tilde{\Gamma}|X}(\gamma|x) dP_{s,u}(\gamma) du$$

$$(2.7) \quad = \int_{-\infty}^v R[f_{\tilde{\Gamma}|X}(\cdot|x)](s, u) du.$$

Here,  $R$  is called the Radon transform (see, *e.g.*, Helgason (1999)), and  $P_{s,u} = \{\gamma : \gamma' s = u\}$  is the affine hyperplane of dimension  $L - 1$  defined through the direction  $s$  and distance  $u$  to the origin, where  $s$  is in  $H^+ = \{x \in \mathbb{S}^{L-1} : x_L > 0\}$  and  $u \in \mathbb{R}$ . The Radon transform is a bounded operator. When applied to a function  $f \in L^1(\mathbb{R}^L)$ , it yields the integral of that function on  $P_{s,u}$

$$R[f](s, u) = \int_{P_{s,u}} f(\gamma) dP_{s,u}(\gamma)$$

where  $dP_{s,u}$  is the Lebesgue measure on  $P_{s,u}$ . Mathematical results regarding this integral transformation (including injectivity and an inversion formula involving the adjoint of the Radon transform) can be found in Natterer (1996) and Helgason (1999). Statistical inverse problems involving this type

of operator on the whole space appear in several problems in tomography (see, *e.g.*, Korostelev and Tsybakov (1993) and Cavalier (2000, 2001)), but also when one wishes to estimate the distribution of random coefficients in the linear model with random coefficients (see Beran, Feuerverger and Hall (1996) and Hoderlein, Klemelä and Mammen (2010)).

**Assumption 2.2.** When  $L = 1$ , for every  $x \in \text{supp}(X)$ ,

$$\text{supp}(f_{V|X}(\cdot|x)) \supset \text{supp}(f_{\Theta|X}(\cdot|x)).$$

When  $L \geq 2$ , for every  $x \in \text{supp}(X)$ ,  $\text{supp}(f_{\tilde{S}|X}(\cdot|x)) = \overline{H^+}$  and for every  $s \in \text{Int}(H^+)$ ,

$$\text{supp}(f_{\tilde{V}|\tilde{S},X}(\cdot|s,x)) \supset \left[ \begin{array}{cc} \inf_{\gamma \in \text{supp}(f_{\tilde{\Gamma}|X}(\cdot|x))} s'\gamma, & \sup_{\gamma \in \text{supp}(f_{\tilde{\Gamma}|X}(\cdot|x))} s'\gamma \end{array} \right].$$

Assumption 2.2 is a large support assumption. It implies that the instruments have a large enough support to apprehend the whole distribution of the unobserved heterogeneity vector. This is crucial in our setup when we want to recover the entire multivariate distribution of heterogeneity factors  $\tilde{\Gamma}$ , and treatment effect parameters for the whole population using weighted averages of conditional the effects UCATE or UCDITE. Assumption 2.2 is not required when  $L = 1$  to obtain treatment effect parameters conditional on the unobserved heterogeneity  $\Theta$ , however, it is required to obtain population averages. When it is not satisfied, we can only make statements about a particular subpopulation related to the variation of the instrument, this is similar to the population apprehended by LATE of Imbens and Angrist (1994), see also Angrist, Grady and Imbens (2000). A similar assumption to Assumption 2.2 is made in Lewbel (2007). Compared to Gautier and Kitamura (2009) and Gautier and Le Penec (2011), Assumption 2.2 allows one of the instruments, for instance  $V$ , to have bounded support.

Finally, Heckman and Smith (1998) consider the case where in a welfare analysis one would like to consider treatment effects in terms of some social welfare criterion  $U$  which could be more general than simply the potential outcomes (*e.g.* consider a more general utility function than income if  $Y$  is income). Within our framework, it is easy to consider the case where the gains are expressed in terms  $U(Y_1, X) - U(Y_0, X)$ , for a known utility function  $U$ , by simply replacing everywhere, including (1.1),  $Y$ ,  $Y_0$ ,  $Y_1$  by  $U(Y, X)$ ,  $U(Y_0, X)$ ,  $U(Y_1, X)$  and transform the variables  $y_i$  into  $U(y_i, x_i)$  for  $i = 1, \dots, N$ .

### 3. IDENTIFICATION OF STRUCTURAL PARAMETERS

This section discusses identification and estimation of the distribution of random coefficients,  $f_{\Gamma', \Theta}$ , the marginal distribution of  $Y_0$  and  $Y_1$ , respectively, given random coefficients,  $f_{Y_j|\Gamma', \Theta}$ ,  $j = 0, 1$ , and several implied parameters, like UCATE. Moreover, we show that the joint distribution and the variance of treatment effects are partially identified, and provide sharp bounds. Finally, we propose an additional assumption that allows us to point identify the variance and distribution of treatment effects, and argue that it is likely to be satisfied in many economically relevant cases.

**3.1. A Central Result for Identification.** We start again by clarifying the notation. In what follows we denote by  $\bar{g}$  the extension of a function  $g$  as 0 outside its domain of definition, *e.g.* a regression function where regressors have bounded support outside of this support. Moreover, denote by  $R$  the Radon transform, and by  $R^{-1}$  its inverse. We use these objects in the following argument, which is at the core of our identification strategy.

First, note that Assumptions (A-1) and (A-2) yield that for  $(v, s', x')$  in  $\text{supp}(\tilde{V}, \tilde{S}', X')$ ,

$$\begin{aligned} \mathbb{E} \left[ \phi(Y)D | \tilde{S} = s, \tilde{V} = v, X = x \right] &= \int_{\text{supp}(\tilde{\Gamma})} \mathbb{E}[\phi(Y_1)\mathbf{1}\{\gamma's < v\} | \tilde{S} = s, \tilde{V} = v, \tilde{\Gamma} = \gamma, X = x] f_{\tilde{\Gamma}|X}(\gamma|x) d\gamma \\ &= \int_{\text{supp}(\tilde{\Gamma})} \mathbf{1}\{\gamma's < v\} \mathbb{E}[\phi(Y_1) | \tilde{S} = s, \tilde{V} = v, \tilde{\Gamma} = \gamma, X = x] f_{\tilde{\Gamma}|X}(\gamma|x) d\gamma \\ &= \int_{\text{supp}(\tilde{\Gamma})} \mathbf{1}\{\gamma's < v\} \mathbb{E}[\phi(Y_1) | \tilde{\Gamma} = \gamma, X = x] f_{\tilde{\Gamma}|X}(\gamma|x) d\gamma. \end{aligned}$$

for any function  $\phi$  such that  $\mathbb{E}[|\phi(Y_1)|] < \infty$ . Thus, by arguments from the previous section,

$$(3.1) \quad \mathbb{E}[\phi(Y)D | \tilde{S} = s, \tilde{V} = v, X = x] = \int_{-\infty}^v R \left[ \overline{\mathbb{E}[\phi(Y_1) | \tilde{\Gamma} = \cdot, X = x] f_{\tilde{\Gamma}|X}(\cdot|x)} \right] (s, u) du.$$

which yields, almost everywhere (a.e. for short) for  $u$  in  $\text{supp}(f_{\tilde{V}|\tilde{S}, X}(\cdot, s, x))$ ,

$$(3.2) \quad \partial_v \mathbb{E}[\phi(Y)D | \tilde{S} = s, \tilde{V} = \cdot, X = x](u) = R \left[ \overline{\mathbb{E}[\phi(Y_1) | \tilde{\Gamma} = \cdot, X = x] f_{\tilde{\Gamma}|X}(\cdot|x)} \right] (s, u).$$

The right hand-side of (3.2) is 0 for  $u$  outside  $\left[ \inf_{\gamma \in \text{supp}(f_{\tilde{\Gamma}|X}(\cdot|x))} s'\gamma, \sup_{\gamma \in \text{supp}(f_{\tilde{\Gamma}|X}(\cdot|x))} s'\gamma \right]$  (see Figure 2). Under Assumption 2.2 we hence know that extending the left hand-side of (3.2) as 0 is innocuous.

These arguments motivate our first theorem. To avoid tedious repetitions, we mostly display results for  $L \geq 2$ . The case of a scalar random coefficient is analogous. It can be obtained by leaving out the inverse of the Radon transform, and adapting the conditioning set accordingly replacing  $\tilde{V}$  and  $\tilde{\Gamma}$  by  $V$  and  $\Theta$ .

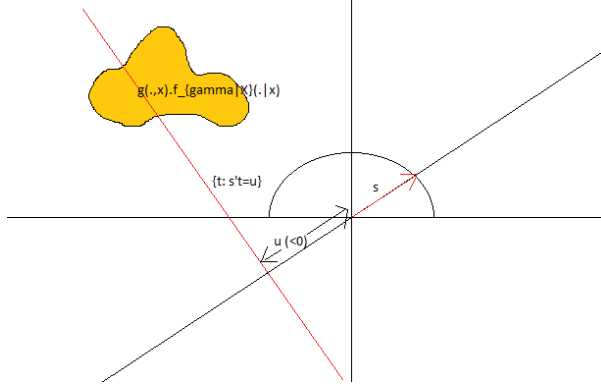


FIGURE 2. Radon transform and extensions

**Theorem 3.1.** Consider an arbitrary function  $\phi$  such that  $\mathbb{E}[|\phi(Y_0)| + |\phi(Y_1)|] < \infty$ . Let  $L \geq 2$ , and assume that Assumptions 2.1 and 2.2 hold. Then, almost surely in  $x$  in  $\text{supp}(X)$ , the following statements are true:

$$(3.3) \quad f_{\tilde{\Gamma}|X}(\cdot|x) = R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ D \mid (\tilde{S}, \tilde{V}) = \cdot, X = x \right]} \right]$$

$$(3.4) \quad \overline{\mathbb{E} \left[ \phi(Y_1) \mid \tilde{\Gamma} = \cdot, X = x \right]} f_{\tilde{\Gamma}|X}(\cdot|x) = R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ \phi(Y) D \mid (\tilde{S}, \tilde{V}) = \cdot, X = x \right]} \right]$$

$$(3.5) \quad \overline{\mathbb{E} \left[ \phi(Y_0) \mid \tilde{\Gamma} = \cdot, X = x \right]} f_{\tilde{\Gamma}|X}(\cdot|x) = R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ \phi(Y) (D - 1) \mid (\tilde{S}, \tilde{V}) = \cdot, X = x \right]} \right].$$

**Discussion of Theorem 3.1:** 1. This set of results always equate a structural parameter of interest on the left hand-side to an object that can be estimated from data on the right hand-side. It is instructive for this first result to compare them with the corresponding results that would be obtained in the scalar unobservable case ( $L = 1$ ),

$$(3.6) \quad f_{\Theta|X}(\cdot|x) = \overline{\partial_v \mathbb{E} [D | V = \cdot, X = x]}$$

$$(3.7) \quad \overline{\mathbb{E} [\phi(Y_1) | \Theta = \cdot, X = x]} f_{\Theta|X}(\cdot|x) = \overline{\partial_v \mathbb{E} [\phi(Y) D | V = \cdot, X = x]}$$

$$(3.8) \quad \overline{\mathbb{E} [\phi(Y_0) | \Theta = \cdot, X = x]} f_{\Theta|X}(\cdot|x) = \overline{\partial_v \mathbb{E} [\phi(Y) (D - 1) | V = \cdot, X = x]}.$$

All the results in the scalar unobservable case involve only one unbounded operator: a partial derivative with respect to  $v$ . In contrast, the results in the multiple unobservable case involve in addition a second unbounded operator: the inverse of the Radon transform, thus showing the more complex ill-posed inverse nature of estimation in this setup. More specifically, in the single unobservable case, the density of the scalar random coefficient  $\Theta$ , conditional on exogenous factors  $X$ , is identified by the derivative of the propensity score with respect to  $v$ , see equation (3.6). In contrast, in order to recover



the (conditional) density of  $\tilde{\Gamma}$  in the case of multiple unobservables, one has to also apply the inverse Radon transform to recover a similar object, see equation (3.3). The same remark applies to the comparison between equations (3.7) and (3.8) in the single unobservable case, and their counterparts (equations (3.4) and (3.5) respectively) in the multiple unobservable case: it is always the inverse Radon transform we have to apply. By trivial manipulations, using the identity for  $\phi$ , we can recover a Heckman and Vytlacil (1998, 2005, 2007) type result,

$$\overline{\mathbb{E}[\Delta|\Theta = \cdot, X = x]} = \frac{\overline{\partial_v \mathbb{E}[Y|V = \cdot, X = x]}}{\overline{\partial_v \mathbb{E}[D|V = \cdot, X = x]}}$$

and can obviously provide an analog in the case of several unobservables. Because of its paramount importance, we focus on the discussion of this conditional average treatment effect in a separate section below.

2. Another important point to notice in Theorem 3.1 is the wide range of functions  $\phi$  that can be used. Using for example  $\phi(y) = \mathbf{1}\{(-\infty, y]\}$  allows to obtain (partial) CDFs, while using  $\phi(t) = \exp(ity)$  allows to obtain partial Fourier transforms, which can be employed to recover densities. This is illustrated in the following corollary.

**Corollary 3.1.** Let Assumptions 2.1 and 2.2 be true, the marginal distributions of  $(Y_0, \tilde{\Gamma})$  and  $(Y_1, \tilde{\Gamma})$  given  $X = x$  are identified. Integrating out  $\tilde{\Gamma}$ , we obtain

$$F_{Y_1|X}(y|x) = \int_{\mathbb{R}^L} R^{-1} \left[ \overline{\partial_v \mathbb{E}[\mathbf{1}\{Y \leq y\} D | (\tilde{S}, \tilde{V}) = \cdot, X = x]} \right] (\gamma) d\gamma,$$

and analogously for  $F_{Y_0|X}(y|x)$ .

From these CDFs we can obtain any inequality measure (*e.g.* the Gini index) for the outcome in the treated and control group. It is important to notice that these quantities are obtained without an ideal randomized experiment. As one example, we may obtain the quantile treatment effect (QTE) of Abadie, Angrist and Imbens (2002), and Chernozhukov and Hansen (2005), which is defined as,

$$\text{QTE}(x, \tau) = q(1, x, \tau) - q(0, x, \tau)$$

where  $q(1, x, \tau)$  and  $q(0, x, \tau)$  are the quantiles of  $F_{Y_1|X}(y|x)$  and  $F_{Y_0|X}(y|x)$ , as well as the related average effect  $\int_0^1 \text{QTE}(x, \tau) d\tau$ , by simply inverting the CDFs obtained from Corollary 3.1.

3. When  $L = 1$ , because we do not have  $R^{-1}$  in the formulas, Assumption 2.2 is not necessary and it is possible to consider cases where the support of the instrument  $V$  is not rich enough to provide identification for every value of the unobservable.

**Proposition 3.1.** Let Assumption 2.1 hold. For almost every  $x$  in  $\text{supp}(X)$  such that  $\mathbb{P}(\Theta \in \text{supp}(f_{V|X}(\cdot|x)) | X = x) > 0$ , we obtain that for every function  $\phi$  s. th.  $\mathbb{E}[|\phi(Y_0)| + |\phi(Y_1)|] < \infty$ ,

$$(3.9) \quad f_{\Theta|X, \Theta \in \text{supp}(f_{V|X}(\cdot|x))}(\cdot|x) = \frac{\overline{\partial_v \mathbb{E}[D|(V = \cdot, X = x)]}}{\int_{\mathbb{R}} \overline{\partial_v \mathbb{E}[D|(V = \cdot, X = x)](t) dt}}$$

$$(3.10) \quad \overline{\mathbb{E}[\phi(Y_1)|\Theta = \cdot, X = x]} f_{\Theta|X, \Theta \in \text{supp}(f_{V|X}(\cdot|x))}(\cdot|x) = \frac{\overline{\partial_v \mathbb{E}[\phi(Y)D|V = \cdot, X = x](t)}}{\int_{\mathbb{R}} \overline{\partial_v \mathbb{E}[D|V = \cdot, X = x](t) dt}},$$

and analogously for  $\overline{\mathbb{E}[\phi(Y_0)|\Theta = \cdot, X = x]}$ . This result allows to identify a large variety of parameters of interest, starting with marginals of potential outcomes

$$F_{Y_0|\Theta \in \text{supp}(f_{V|X}(\cdot|x)), X}(y|x) = \frac{\int_{\text{supp}(f_{V|X}(\cdot|x))} \partial_v \mathbb{E}[\mathbf{1}\{Y \leq y\}(D-1)|V = \cdot, X = x](t) dt}{\int_{\text{supp}(f_{V|X}(\cdot|x))} \partial_v \mathbb{E}[D|V = \cdot, X = x](t) dt}.$$

and including average or quantile treatment effects, variance and distribution of treatment effects (under the respective assumptions to point identify these quantities that we make below). The subpopulation for which we can make inference is the same as in Angrist, Grady and Imbens (2000). We conjecture that an analog result holds in the case of a multivariate unobservable.

To close this discussion, we remark that the equalities in Theorem 3.1 hold a.e., for convenience we will no longer mention this in the remainder of the paper. In the following, we discuss a key implication of Theorem 3.1. Due to its great importance, we state it in a separate subsection.

**3.2. The Average Treatment Effect Conditional on Unobservables (UCATE).** In this section we extend the notion of the MTE in Björklund and Moffitt (1987) and Heckman and Vytlacil (1998, 2005, 2007) to our setup with potentially several unobservables. In the presence of only one unobservable in the selection equation, we call the parameter

$$\text{UCATE}(\theta, x) := \mathbb{E}[\Delta | \Theta = \theta, X = x],$$

the Unobservables Conditioned Average Treatment Effect (UCATE, for short). It is the average effect of treatment for the subpopulation with unobserved heterogeneity parameter equal to  $\theta$  and observable  $X = x$ . The parameter  $\text{UCATE}(\theta, x)$  has the same interpretation as the MTE. It corresponds to the average effect for a subpopulation with  $X = x$  who would be indifferent between participation and non-participation in the treatment, if they were exogenously assigned a value  $v$  of  $V$  such that  $v = \theta$ . The UCATE parameter has the advantage that it can be generalized easily to the case where  $L \geq 2$  as

$$\text{UCATE}(\gamma, x) := \mathbb{E}[\Delta | \tilde{\Gamma} = \gamma, X = x].$$

It is the average effect of treatment for the subpopulation with first stage unobserved heterogeneity vector equal to  $\gamma$  and observables  $X = x$ . It is also the average effect for people with  $X = x$  who would be indifferent between participation and non-participation in the treatment, if they were exogenously assigned a value  $(s, v)$  of  $(\tilde{S}, \tilde{V})$ , such that  $s'\gamma = v$ . Because Assumption (A-2) yields that

$$\begin{aligned} & \mathbb{E} \left[ \Delta | \tilde{\Gamma} = \gamma, X = x, (\tilde{S}, \tilde{V}) = (s, v) \right] \\ &= \mathbb{E} \left[ Y_1 | \tilde{\Gamma} = \gamma, X = x, (\tilde{S}, \tilde{V}) = (s, v) \right] - \mathbb{E} \left[ Y_0 | \tilde{\Gamma} = \gamma, X = x, (\tilde{S}, \tilde{V}) = (s, v) \right] \\ &= \mathbb{E}[Y_1 | \tilde{\Gamma} = \gamma, X = x] - \mathbb{E}[Y_0 | \tilde{\Gamma} = \gamma, X = x] \\ &= \mathbb{E}[\Delta | \tilde{\Gamma} = \gamma, X = x], \end{aligned}$$

UCATE( $\gamma, x$ ) does not depend on the values taken by the instruments  $(\tilde{S}, \tilde{V})$ . It is thus policy invariant. All these properties are essential properties of the MTE. The extensions of MTE suggested in Heckman and Vytlacil (2005) for non additively separable models in the first stage selection equation do not satisfy these properties simultaneously so we believe that UCATE is a natural parameter to extend the standard treatment effects analysis to more general models for the treatment choice. Moreover, like in Heckman and Vytlacil (2005), a large variety of treatment effect measures that depend on averages can be written as weighted averages of the UCATE parameter. The ATE for example can be obtained from UCATE under Assumption 2.2. To this end, we employ the general results in the previous subsection, with  $\phi(y) = y$ ,  $\forall y \in \mathbb{R}$  and make use of the following assumption.

**Assumption 3.1.**  $\mathbb{E}[|Y_1| + |Y_0|] < \infty$ .

In the following, we derive UCATE formally, which is straightforward given the results above.

**Theorem 3.2.** Suppose assumptions 2.1 and 3.1 hold. When  $L = 1$ , we obtain for almost every  $\theta$ ,

$$(3.11) \quad \overline{\text{UCATE}(\theta, x)} f_{\Theta|X}(\theta|x) = \partial_v \mathbb{E} [Y | V = \cdot, X = x] (\theta).$$

Suppose now that  $L \geq 2$  and invoke Assumption 2.2. Then, for almost every  $\gamma$  in  $\mathbb{R}^L$ ,

$$(3.12) \quad \overline{\text{UCATE}(\gamma, x)} f_{\tilde{\Gamma}|X}(\gamma|x) = R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ Y \mid (\tilde{S}, \tilde{V}) = \cdot, X = x \right]} \right] (\gamma).$$

Obviously, from this equation it is trivial to solve for UCATE by simply plugging in the expressions for  $f_{\Theta|X}(\gamma|x)$  (respectively  $f_{\tilde{\Gamma}|X}(\gamma|x)$ ) from Theorem 3.1 into this equation. The right hand-side of (3.12) is a natural extension of the local instrumental variable estimator (LIV for short). If UCATE( $\gamma, x$ ) is found to be not constant in one dimension of  $\gamma$ , it is an indication of heterogeneous

“costs factors” in this specific direction. We then conclude that accounting for it through a random coefficients specification is essential for recovery of the classical treatment effect parameters.

UCATE is a building block from which we may derive a large variety of treatment effect parameters that depend on averages, *e.g.*, the average treatment effect, or the treatment effect on the treated:

$$\begin{aligned} \text{ATE}(x) &= \int_{\mathbb{R}^L} \overline{\text{UCATE}(\gamma, x)} f_{\tilde{\Gamma}|X}(\gamma|x) d\gamma = \int_{\mathbb{R}^L} R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ Y \mid (\tilde{S}, \tilde{V}) = \cdot, X = x \right]} \right] (\gamma) d\gamma \\ \text{TT}(x) &= \int_{\mathbb{R}^L} h_{\text{TT}}(\gamma, x) \overline{\text{UCATE}(\gamma, x)} f_{\tilde{\Gamma}|X}(\gamma|x) d\gamma = \int_{\mathbb{R}^L} h_{\text{TT}}(\gamma, x) R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ Y \mid (\tilde{S}, \tilde{V}) = \cdot, X = x \right]} \right] (\gamma) d\gamma \end{aligned}$$

where  $h_{\text{TT}}(\gamma, x) = \mathbb{E} \left[ \mathbf{1} \left\{ \tilde{S}'\tilde{\Gamma} < \tilde{V} \right\} \mid X = x \right]^{-1} \mathbb{E} \left[ \mathbf{1} \left\{ \tilde{S}'\gamma < \tilde{V} \right\} \mid X = x \right]$ , and analogously for the treatment effect on the untreated. In Section 3.4, we discuss a natural extension of UCATE to obtain more general treatment effect parameters which depend on the whole distribution of treatment effects, or on the full joint distribution of  $(Y_0, Y_1, \tilde{\Gamma})$ .

This concludes our discussion of the point identified effects; however, we can use the results in Theorem 3.1 in additional ways, as the following subsection illustrates.

**3.3. Bounds on the Variance and the Distribution of Treatment Effects.** As discussed above, Theorem 3.1 reveals that under our assumptions, the marginals of  $Y_1$  and  $Y_0$ , conditional on preference parameters, are identified. This information about the marginals, as well as moments of the marginals, can be used to bound moments of the difference  $Y_1 - Y_0$ , as well as the distribution of this difference, as this section illustrates.

For the formal statement of the bounds, we require the following notation. Let

$$\begin{aligned} F^L(y_0, y_1|x) &= \max \left\{ \int_{\mathbb{R}^L} R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ \mathbf{1} \{Y \leq y\} (2D - 1) \mid (\tilde{S}, \tilde{V}) = \cdot, X = x \right]} \right] (\gamma) d\gamma - 1, 0 \right\} \\ F^U(y_0, y_1|x) &= \min \left\{ \int_{\mathbb{R}^L} R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ \mathbf{1} \{Y \leq y\} (D - 1) \mid (\tilde{S}, \tilde{V}) = \cdot, X = x \right]} \right] (\gamma) d\gamma, \right. \\ &\quad \left. \int_{\mathbb{R}^L} R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ \mathbf{1} \{Y \leq y\} D \mid (\tilde{S}, \tilde{V}) = \cdot, X = x \right]} \right] (\gamma) d\gamma \right\} \end{aligned}$$

and

$$\begin{aligned} F^L(\delta|x) &= \int_{\mathbb{R}^L} \sup_{y \in \mathbb{R}} \max \left\{ R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ \mathbf{1} \{Y \leq y\} D \mid (\tilde{S}, \tilde{V}) = \cdot, X = x \right]} \right] (\gamma) \right. \\ &\quad \left. - R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ \mathbf{1} \{Y \leq y - \delta\} (D - 1) \mid (\tilde{S}, \tilde{V}) = \cdot, X = x \right]} \right] (\gamma), 0 \right\} d\gamma \end{aligned}$$

$$F^U(\delta|x) = 1 + \int_{\mathbb{R}^L} \inf_{y \in \mathbb{R}} \min \left\{ R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ \mathbf{1} \{Y \leq y\} D \mid (\tilde{S}, \tilde{V}) = \cdot, X = x \right]} \right] (\gamma) \right.$$

$$-R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ \mathbf{1} \{Y \leq y - \delta\} (D - 1) \mid \left( \tilde{S}, \tilde{V} \right) = \cdot, X = x \right]} \right] (\gamma), 0 \Big\} d\gamma.$$

Now we are in a position to characterize the bounds for the variance and the distribution of treatment effects in our baseline scenario.

**Theorem 3.3.** Suppose that Assumptions 2.1 and 2.2 hold. Then we obtain that, almost surely in  $x$  in  $\text{supp}(X)$ , for every  $(y_0, y_1, \delta) \in \mathbb{R}^3$ ,

$$(3.13) \quad F^L(y_0, y_1 | x) \leq F_{Y_0, Y_1 | X}(y_0, y_1 | x) \leq F^U(y_0, y_1 | x)$$

and

$$(3.14) \quad F^L(\delta | x) \leq F_{\Delta | X}(\delta | x) \leq F^U(\delta | x),$$

if in addition  $\mathbb{E}[(Y_0)^2 + (Y_1)^2] < \infty$ , then

$$\begin{aligned} & \left( \text{Var}(\Delta | X = x) + \left( \int_{\mathbb{R}^L} R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ Y \mid \left( \tilde{S}, \tilde{V} \right) = \cdot, X = x \right]} \right] (\gamma) d\gamma \right)^2 \right. \\ & + \left. \int_{\mathbb{R}^L} R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ (1 - 2D)Y \mid \left( \tilde{S}, \tilde{V} \right) = \cdot, X = x \right]} \right] (\gamma) d\gamma \right)^2 \\ & \leq 4 \int_{\mathbb{R}^L} R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ (D - 1)Y^2 \mid \left( \tilde{S}, \tilde{V} \right) = \cdot, X = x \right]} \right] (\gamma) d\gamma \int_{\mathbb{R}^L} R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ DY^2 \mid \left( \tilde{S}, \tilde{V} \right) = \cdot, X = x \right]} \right] (\gamma) d\gamma. \end{aligned}$$

Note that this result implies that we can bound the variance and the distribution of treatment effects,  $\text{Var}(\Delta | X = x)$  and  $F_{\Delta | X}(\delta | x)$ , respectively, entirely by observable quantities. The bounds on the variance are very easy to obtain and do not need to be obtained from the CDFs. The bound (3.13) is a direct application of the Fréchet-Hoeffding bounds, like Heckman, Smith and Clements (1997), Manski (1997) and Heckman and Smith (1998). It is obtained from the conditional CDFs  $F_{Y_0 | X}(y_0, x) = \mathbb{E}_{\tilde{\Gamma}} \left[ F_{Y_0 | X, \tilde{\Gamma}}(y_0 | x, \tilde{\Gamma}) \right]$  and  $F_{Y_1 | X}(y_1 | x) = \mathbb{E}_{\tilde{\Gamma}} \left[ F_{Y_1 | X, \tilde{\Gamma}}(y_1 | x, \tilde{\Gamma}) \right]$ . The bound (3.14) is, like Fan and Park (2010) and Firpo and Ridder (2008), a consequence of the Makarov bounds. Firpo and Ridder (2008) show that, unlike the Fréchet-Hoeffding bounds which are uniformly sharp, the Makarov bounds are only pointwise sharp (unless the outcomes are binary). We show that the average Makarov bounds are pointwise sharp and tighter than the Makarov bounds on the average distribution. (3.14) is obtained by taking the expectation of the Makarov bounds evaluated at  $F_{\Delta | X, \tilde{\Gamma}}(\delta | x, \tilde{\Gamma})$ . In a similar fashion, we derive sharper bounds on  $F_{\Delta}$  from (3.14) by averaging over covariates than what we would obtain by calculating the Makarov bounds from the CDFs of  $Y_0$  and  $Y_1$ , *i.e.*,

$$\mathbb{E}_X [F^L(\delta | X)] \leq F_{\Delta}(\delta) \leq \mathbb{E}_X [F^U(\delta | X)].$$

(3.13) and (3.14) are obtained integrating out the unobservables, they take the random coefficients structure in the selection equation into account. While the above mentioned references present bounds in the case of randomized experiments or selection on observables, we obtain for the first time bounds in the case where there could be endogenous selection into treatment.

Confidence bands and bounds on functionals of  $F_{\Delta|X}$  can be obtained in a similar spirit as Fan and Park (2010) and Firpo and Ridder (2008). Bounds on functionals of  $F_{Y_0, Y_1|X}$  could be obtained in similar spirit of Fan and Zhu (2009). Unlike these references, these bounds would hold in the case of endogenous selection. They would also be tighter because we exploit the specific random coefficients structure in the selection equation. Nevertheless, these bounds could be large in practice. Hence, we now discuss possible alternatives that provide point identification.

**3.4. Point Identification of Variance and Distribution of Treatment Effects.** In this section, as well as in Sections 4.3.3 and 4.3.4, we consider the particular situation where the random variables  $Y_0$  and  $Y_1$  are continuously distributed. We thus replace (A-1) by (A-(1)') in Assumption 2.1.

As we have just seen, in order to point identify the variance and the distribution of treatment effects, it is imperative to invoke further assumptions. In our endogenous setup, these assumptions involve the unobservables. For the identification of the distribution of treatment effects, we make the following key assumption.

**Assumption 3.2.**  $Y_0 \perp \Delta \mid \tilde{\Gamma}, X$

A slightly weaker form is sufficient for point identification of the variance of treatment effects. Heckman, Smith and Clements (1997) considers independence of the base state and the gains given  $D$ . We argue below that controlling for a vector of random coefficients instead of  $D$  retains the same spirit as in Heckman, Smith and Clements (1997), but makes the assumption more plausible. Heckman and Smith (1998) consider independence given a vector of observable characteristics  $X$ . As we argue subsequently in various extensions of the Roy model, it is important to control for the unobserved heterogeneity that enter in the selection equation as well. To keep the notation minimal, we suppress henceforth the dependence on  $X$ .

Assumption 3.2 can be readily interpreted in terms of control functions. It is useful to think of  $Y$  as being generated by a nonseparable model; in this case  $Y = \psi(D, U)$ , and  $Y_0 = \psi(0, U)$ , as well as  $Y_1 = \psi(1, U) = Y_0 + \Delta$ . If we identify  $U$  with a high dimensional unobservable, *e.g.*, preferences, it is interesting to note that this implies that our random coefficients  $Y_0$  and  $\Delta$  are two different functions of these unobservables, *i.e.*,  $Y_0 = a(U) = \psi(0, U)$  and  $\Delta = b(U) = \psi(1, U) - \psi(0, U)$ . Without loss

of generality, one could further partition the set of unobservable in vectors  $U_0, U_1$ , and  $U_2$ , and write  $Y_0 = a(U_0, U_2)$  and  $\Delta = b(U_1, U_2)$ .

Assumption 3.2 restricts the heterogeneity appearing in this model. It is best understood in terms of the reformulation introduced above, namely  $Y_0 = a(U_0, U_2)$  and  $\Delta = b(U_1, U_2)$ . A sufficient condition for Assumption (3.2) is that  $(\Gamma', \Theta) = U_2$ , and  $U_0 \perp U_1 | (\Gamma', \Theta)$  (for the latter it would in turn be sufficient that  $U_0 \perp U_1 \perp (\Gamma', \Theta)$ ). In words, there is a common driving factor that causes the selection bias and the dependence between  $Y_0$  and  $\Delta$ . This factor is given by  $(\Gamma', \Theta)$ , which, even though it is not recovered for every individual, implicitly serves as a control function. There is remaining randomness in  $Y_0$  and  $\Delta$ , however, once the driving factor for endogeneity in this system, *i.e.*,  $(\Gamma', \Theta)$ , is accounted for, there is no leftover dependence.

Note that it does **not** mean that  $Y_0 \perp \Delta$ . In fact, unless there is no endogenous selection there will generally be dependence between  $Y_0$ , and  $Y_1 - Y_0$ . In summary, there is endogenous selection into treatment, but as far as it is endogenous, it can be summarized by  $(\Gamma', \Theta)$ . Note that the assumption is more likely to hold in the model with several unobservables in the selection equation, in the sense that there is not just a single factor that we can employ to control for endogeneity, but a full vector of such variables.

The next subsection details that this assumption is sensible in several extensions of the Roy model.

3.4.1. *The Example of The Roy Model.* The aim of this section is to show that Assumption 3.2 is satisfied for several extensions of the popular Roy model. Consider a model where the individuals have at their disposal an information set  $\mathcal{I}$  at the time of their decision to participate in the program. They select themselves into treatment if and only if their expected net utility exceeds expected costs.

Formally, the individuals choose  $D = 1$  if and only if

$$\mathbb{E}[Y_1 - Y_0 - C_1 | \mathcal{I}] > 0,$$

where  $C_1$  denotes the costs associated with participating in treatment. For simplicity, throughout this subsection, we do not assume to condition on observed factors  $X$ . However, we note that they could be used to make Assumption 3.2 more plausible. We use the notation  $\Delta = \mathbb{E}[\Delta | \mathcal{I}] + \Xi$ . We also assume that the expected costs  $\mathbb{E}[C_1 | \mathcal{I}]$  can be approximated on individual level by a linear function, *i.e.*,  $\mathbb{E}[C_1 | \mathcal{I}] \cong \bar{\Gamma}_0 - \bar{\Gamma}_1 V + \bar{\Gamma}' Z$ , where  $\bar{\Gamma}_1 > 0$  almost surely (the original  $V$  can be changed to  $-V$ ). Since the population is heterogeneous, these coefficients vary across the population. Dividing the

expected net utility of treatment by  $\bar{\Gamma}_1$ , we obtain the selection equation (1.2) with  $\Gamma = \bar{\Gamma}/\bar{\Gamma}_1$  and  $\Theta = (\bar{\Gamma}_0 - \mathbb{E}[\Delta|\mathcal{I}])/\bar{\Gamma}_1$ .

Let us consider a first setup where Assumption 3.2 is satisfied. Suppose that  $\bar{\Gamma}_0 - \bar{\Gamma}_1 V + \bar{\Gamma}' Z = h_0(\Psi) + h_1(\Psi)V + h_2(\Psi)'Z$  where  $\Psi$  denotes some deep economic parameters. When  $L \geq 2$  we have

$$\Gamma = \frac{h_2(\Psi)}{h_1(\Psi)} = \left( \frac{h_{2,1}(\Psi)}{h_1(\Psi)}, \dots, \frac{h_{2,L-1}(\Psi)}{h_1(\Psi)} \right)'.$$

It is reasonable to believe that when  $L$  is relatively large and the instruments are well chosen  $\Gamma$  captures a lot of features of the deep parameters  $\Psi$ . The following assumption considers an ideal situation.

**Assumption 3.3** (Invertibility). There is a bijective mapping from  $\Psi$  into  $\Gamma = h_2(\Psi)/h_1(\Psi)$ .

It implies that  $\Psi$  and hence also  $\bar{\Gamma}_0$  and  $\bar{\Gamma}_1$  are  $\sigma(\Gamma)$ -measurable, where  $\sigma(A)$  denotes the sigma algebra generated by a random vector  $A$ . Hence  $\mathbb{E}[\Delta|\mathcal{I}] = -\bar{\Gamma}_1\Theta + \bar{\Gamma}_0$  is  $\sigma(\Gamma, \Theta)$ -measurable. Note that we do not have to know or estimate this mapping.

Assume as well

**Assumption 3.4.**  $\Xi$  is  $\sigma(\Gamma, \Theta)$ -measurable.

When  $\mathcal{I} \supset \sigma(\Gamma, \Theta)$  Assumption 3.4 can be rewritten in the form  $\Xi = 0$ , *i.e.* the agents have perfect foresight.

**Proposition 3.2.** Assumptions 3.3 and 3.4 imply Assumption 3.2. Also, Assumption 3.2 is satisfied even if we switch the labels between 0 and 1.

Proposition 3.2 is a direct consequence of the decomposition  $\Delta = \mathbb{E}[\Delta|\mathcal{I}] + \Xi$ . Indeed the assumptions yield that  $\Delta$  is  $\sigma(\Gamma, \Theta)$ -measurable and  $Y_0 \perp \Delta|\Gamma, \Theta$  is trivially satisfied as conditional on the unobservables the treatment effect is constant. Having an assumption that is independent of the labeling is a desirable property when there is no specific treatment but simply two different states, *e.g.*, two different employment sectors.

In the second setup we assume more generally that the forecast error on the gains is independent of the outcome in base state, given the rescaled sources of unobserved heterogeneity in the selection equation.

**Assumption 3.5.**  $\Xi \perp Y_0|\Gamma, \Theta$ .

The following proposition simply relies on the decomposition  $\Delta = \mathbb{E}[\Delta|\mathcal{I}] + \Xi$ .



**Proposition 3.3.** Assumptions 3.3 and 3.5 imply Assumption 3.2.

Note, however, that if  $\Gamma$  and  $\Theta$  are known to the agents at the time the decision is made,  $\mathbb{E}[\Delta|Z] = \text{UCATE}$  is identified and  $f_{\Xi|\tilde{\Gamma}}$  is identified if  $f_{\Delta|\tilde{\Gamma}}$  is identified, see Section 3.4.3. Invertibility has thus a strong implication regarding the structure of the ex-ante information set.

In the third setup we relax Assumption 3.3 from setups 1 and 2 to allow  $\bar{\Gamma}_0, \bar{\Gamma}_1$  to not be  $\sigma(\Gamma)$ -measurable. The following proposition gives a more general sufficient condition for Assumption 3.2 to hold that is satisfied under Assumptions 3.3 and 3.5.

**Proposition 3.4.** Assume that  $\bar{\Gamma}_1 = \lambda(\Gamma, \Xi_1)$  for some measurable function  $\lambda$ , and that  $(\bar{\Gamma}_0, \Xi, \Xi_1) \perp Y_0|\Gamma, \Theta$ , then Assumption 3.2 is satisfied.

The proposition follows from the fact that

$$\begin{aligned} \Delta &= \mathbb{E}[\Delta|Z] + \Xi \\ &= -\lambda(\Gamma, \Xi_1)\Theta + \bar{\Gamma}_0 + \Xi. \end{aligned}$$

Note that, when in the original scale one coordinate of  $\bar{\Gamma}$  is non random, then  $\bar{\Gamma}_1 \in \sigma(\Gamma)$  and thus there exists a measurable function  $\lambda$  such that  $\bar{\Gamma}_1 = \lambda(\Gamma, \Xi_1)$ . Proposition (3.4) is also satisfied when  $\bar{\Gamma}_1$  is non random or when  $\bar{\Gamma}_1 = \Xi_1$  and  $(\bar{\Gamma}_0, \Xi, \Xi_1) \perp Y_0|\Gamma, \Theta$ . Moreover, it is worth mentioning that there could be arbitrary dependence between  $(\bar{\Gamma}_0, \Xi, \Xi_1)$  holding fixed  $\Gamma, \Theta$  and that we do not assume that  $(\bar{\Gamma}_0, \Xi, \Xi_1) \perp Y_0$  but rather that they can only depend on each other through  $\Gamma, \Theta$ .

The assumptions of Proposition 3.4 allow for more complex structure of the ex-ante information set than Assumptions 3.3 and 3.5. Indeed  $\mathbb{E}[\Delta|Z] = -\lambda_1(\Gamma, \Xi_1)\Theta + \bar{\Gamma}_0$  can be different from UCATE if  $\lambda$  is non constant in its last argument and/or constant in others (*e.g.* the agent does not fully know some of his  $\Gamma$ 's).

Proposition 3.4 is the most general sufficient condition for Assumption 3.2 to hold that we present. It could certainly be further generalized. Thus there is a wide class of structural models that encompass generalizations of the Roy model for which Assumption 3.2 can hold.

**3.4.2. Variance of the Treatment Effect.** In the following we show that the Unobservables Conditioned Variance of the treatment effects, called UCVATE and defined as  $Var(\Delta|\Theta = \theta, X = x)$  and  $Var(\Delta|\tilde{\Gamma} = \gamma, X = x)$ , in the cases of  $L = 1$ , and  $L \geq 2$ , respectively, is point identified under an assumption that is similar in spirit, but weaker than the full independence assumption specified above. It is the following conditional uncorrelatedness.

**Assumption 3.6.**  $\mathbb{E}[Y_0^2 + Y_1^2] < \infty$ ,  $\mathbb{E}[Y_0\Delta | \Gamma, \Theta, X] = \mathbb{E}[Y_0 | \Gamma, \Theta, X] \mathbb{E}[\Delta | \Gamma, \Theta, X]$ . When  $L = 1$ , there are no  $\Gamma$ 's in the conditioning set.

For the sake of completeness, we retain the case of a single unobservable in the selection equation.

**Theorem 3.4.** Let Assumptions 2.1 and 3.6 hold, and  $L = 1$ . Then, almost surely in  $x$  in  $\text{supp}(X)$ , for almost every  $\theta \in \text{supp}(\Theta)$

$$\begin{aligned} \text{Var}(\Delta | \Theta = \theta, X = x) f_{\Theta|X}(t|x) &= \overline{\partial_v \mathbb{E}[Y^2 | V = \cdot, X = x]}(t) \\ &+ \frac{\overline{\partial_v \mathbb{E}[Y | V = \cdot, X = x]}(t) \overline{\partial_v \mathbb{E}[(1 - 2D)Y | V = \cdot, X = x]}(t)}{f_{\Theta|X}(t|x)}. \end{aligned}$$

In the case where  $L \geq 2$ , if we also invoke Assumption 2.2, we obtain that for almost every  $\gamma \in \text{supp}(\tilde{\Gamma})$ ,

$$\begin{aligned} \text{Var}(\Delta | \tilde{\Gamma} = \gamma, X = x) f_{\tilde{\Gamma}|X}(\gamma|x) &= R^{-1} \left[ \overline{\partial_v \mathbb{E}[Y^2 | (\tilde{S}, \tilde{V}) = \cdot, X = x]} \right] (\gamma) \\ &+ \frac{R^{-1} \left[ \overline{\partial_v \mathbb{E}[Y | (\tilde{S}, \tilde{V}) = \cdot, X = x]} \right] (\gamma) R^{-1} \left[ \overline{\partial_v \mathbb{E}[(1 - 2D)Y | (\tilde{S}, \tilde{V}) = \cdot, X = x]} \right] (\gamma)}{f_{\tilde{\Gamma}|X}(\gamma|x)}. \end{aligned}$$

Finally, under the same additional Assumption 2.2, it holds that

$$\text{Var}(\Delta | X = x) = \int_{\mathbb{R}} \text{Var}(\Delta | \Theta = \theta, X = x) f_{\Theta|X}(t|x) dt + \int_{\mathbb{R}} (\mathbb{E}[\Delta | \Theta = \theta, X = x] - \text{ATE}(x))^2 f_{\Theta|X}(t|x) d\gamma.$$

when  $L = 1$ , and

$$\text{Var}(\Delta | X = x) = \int_{\mathbb{R}^L} \text{Var}(\Delta | \tilde{\Gamma} = \gamma, X = x) f_{\tilde{\Gamma}|X}(\gamma|x) d\gamma + \int_{\mathbb{R}^L} (\mathbb{E}[\Delta | \tilde{\Gamma} = \gamma, X = x] - \text{ATE}(x))^2 f_{\tilde{\Gamma}|X}(\gamma|x) d\gamma.$$

when  $L \geq 2$ .

These rather involved formulae provide nevertheless a succinct description of the conditional variance, which does not require material assumptions beyond instrument independence, no correlation between base state and treatment effect, conditional on all observable and unobservable variables (Assumption 3.6), and a specific relation between the variation of the instruments and the support of the unobserved heterogeneity in the model for the selection into treatment (Assumption 2.2). When  $L = 1$ , but the latter assumption does not hold, we may again identify the variance of the treatment effect for the population such that  $\Theta \in \text{supp}(f_{V|X}(\cdot|x))$  given  $X = x$ .

3.4.3. *UCDITE and Treatment Effect Parameters that Depend on the Distribution of Treatment Effects or the Joint Distribution of Potential Outcomes.* This section extends the analysis of the previous subsections to distributions of treatment effects. We define the Conditional Distribution of Treatment Effects as

$$\text{UCDITE}(\delta, \theta, x) := f_{\Delta|\Theta, X}(\delta|\theta, x)$$

when  $L = 1$  and

$$\text{UCDITE}(\delta, \gamma, x) := f_{\Delta|\tilde{\Gamma}, X}(\delta|\gamma, x)$$

when  $L \geq 2$ . This quantity is the distribution of treatment effects for the subpopulation with unobserved heterogeneity parameter (respectively vector) from the first stage selection equation equal to  $\theta$  (respectively  $\gamma$ ) and observables  $X = x$ . It is also the distribution of the gains in terms of  $Y_1 - Y_0$  for people with  $X = x$  who would be indifferent between participation and non-participation in the treatment if they were exogenously assigned the value  $v$  of  $V$  (respectively  $(s, v)$  of  $(\tilde{S}, \tilde{V})$ ), such that  $v = \theta$  (respectively  $s'\gamma = v$ ). It is straightforward to check that, akin to UCATE, Assumption (A-2) yields that  $\text{UCDITE}(\theta, x)$  (respectively  $\text{UCDITE}(\gamma, x)$ ) does not depend on the values taken by the instrument  $V$  (respectively  $(\tilde{S}, \tilde{V})$ ), and is hence policy invariant. This quantity is at the heart of any calculation of more general treatment effects parameters that go beyond averages and depend on the distribution of either the treatment effects  $Y_1 - Y_0$ , or the joint of potential outcomes  $(Y_1, Y_0)$ .

To calculate some of these distributional treatment effects we sometimes have to replace Assumption 2.1 (A-2) by the following stronger assumption.

**Assumption 3.7.**

$$(V, Z) \perp (Y_0, Y_1, \Gamma', \Theta) | X.$$

More specifically, out of the list of the following parameters, which can be deduced from UCDITE, Assumption 3.7 is required to obtain the last three,

$$\begin{aligned} f_{\Delta|X}(\delta|x) &= \int_{\mathbb{R}^L} \overline{\text{UCDITE}(\delta, \gamma, x)} f_{\tilde{\Gamma}|X}(\gamma|x) d\gamma \\ \mathbb{P}(\Delta > 0 | X = x) &= \int_{\mathbb{R}} \mathbf{1}\{\delta > 0\} \int_{\mathbb{R}^L} \overline{\text{UCDITE}(\delta, \gamma, x)} f_{\tilde{\Gamma}|X}(\gamma|x) d\gamma d\delta \\ f_{Y_0, Y_1|X}(y_0, y_1|x) &= \int_{\mathbb{R}^L} \overline{\text{UCDITE}(y_1 - y_0, \gamma, x) f_{Y_0|\tilde{\Gamma}, X}(y_0|\gamma, x)} f_{\tilde{\Gamma}|X}(\gamma|x) d\gamma \\ f_{\Delta|D=1, X}(\delta|x) &= \int_{\mathbb{R}^L} h_{\text{TT}}(\gamma, x) \overline{\text{UCDITE}(\delta, \gamma, x)} f_{\tilde{\Gamma}|X}(\gamma|x) d\gamma \\ f_{\Delta|D=1, Y_0, X}(\delta|y_0, x) &= \int_{\mathbb{R}^L} h_{\text{TT}}(\gamma, x) \overline{\text{UCDITE}(\delta, \gamma, x) f_{Y_0|\tilde{\Gamma}, X}(y_0|\gamma, x)} f_{\tilde{\Gamma}|X}(\gamma|x) d\gamma. \end{aligned}$$

At this stage, it is important to recall that  $f_{Y_0|\tilde{\Gamma},X}(y_0|\gamma,x)f_{\tilde{\Gamma}|X}(\gamma|x)$  may be obtained from Theorem 3.1<sup>7</sup>. Setting  $\phi(y) = e^{ity}$  in Theorem 3.1 yields, almost surely in  $x$  in  $\text{supp}(X)$ ,

$$\begin{aligned}\mathcal{F}_1 [f_{Y_0+\Delta,\Theta}] (t, \cdot) &= \overline{\partial_v \mathbb{E} [e^{itY} D | V = \cdot, X = x]} \\ \mathcal{F}_1 [f_{Y_0,\Theta}] (t, \cdot) &= \overline{\partial_v \mathbb{E} [e^{itY} (D - 1) | V = \cdot, X = x]},\end{aligned}$$

while, when  $L \geq 2$ , under Assumption 2.2, we obtain

$$\begin{aligned}\mathcal{F}_1 [f_{Y_0+\Delta,\tilde{\Gamma}|X}(\cdot|x)] (t, \cdot) &= R^{-1} \left[ \overline{\partial_v \mathbb{E} [e^{itY} D | (\tilde{S}, \tilde{V}) = \cdot, X = x]} \right] \\ \mathcal{F}_1 [f_{Y_0,\tilde{\Gamma}|X}(\cdot|x)] (t, \cdot) &= R^{-1} \left[ \overline{\partial_v \mathbb{E} [e^{itY} (D - 1) | (\tilde{S}, \tilde{V}) = \cdot, X = x]} \right].\end{aligned}$$

where  $\mathcal{F}_1$  denotes the Fourier transform of the joint density seen as a function of its first variable, holding the other arguments fixed. This object is called a partial Fourier transform.

The importance of Assumption 3.2 is that it allows factorization, *i.e.*,

$$\mathcal{F}_1 [f_{Y_0+\Delta,\tilde{\Gamma}|X}(\cdot|x)] = \mathcal{F}_1 [f_{Y_0,\tilde{\Gamma}|X}(\cdot|x)] \mathcal{F}_1 [f_{\Delta,\tilde{\Gamma}|X}(\cdot|x)],$$

when  $L \geq 2$ . We make moreover use of the following integrability assumption.

**Assumption 3.8.** For almost every  $x$  in  $\text{supp}(X)$ , for every  $\theta$  in  $\text{supp}(f_{\theta|X}(\cdot|x))$ , respectively for every  $\gamma$  in  $\text{supp}(f_{\tilde{\Gamma}|X}(\cdot|x))$ ,  $\text{UCDITE}(\cdot, \theta, x)$ , respectively  $\text{UCDITE}(\cdot, \gamma, x)$ , belong to  $L^1(\mathbb{R}) \cap L^2(\mathbb{R})$ .

Finally, we require a last technical assumption, which is common in the deconvolution literature.

**Assumption 3.9.** When  $L = 1$ ,

$$\forall x \in \text{supp}(X), \forall (t, \theta) \in \mathbb{R} \times \text{supp}(f_{\theta|X}(\cdot|x)), \mathcal{F} [f_{Y_0|\Theta,X}(\cdot|\theta, x)] (t) \neq 0$$

while when  $L \geq 2$ ,

$$\forall x \in \text{supp}(X), \forall (t, \gamma) \in \mathbb{R} \times \text{supp}(f_{\tilde{\Gamma}|X}(\cdot|x)), \mathcal{F} [f_{Y_0|\tilde{\Gamma},X}(\cdot|\gamma, x)] (t) \neq 0$$

where  $\mathcal{F}$  is the Fourier transform.

We believe that it is possible to weaken this assumption and allow for isolated zeros in the spirit of Devroye (1989), Carrasco and Florens (2011) and Evdokimov and White (2011) among others but prefer not to elaborate on this in this article.

<sup>7</sup>Similar expressions can be obtained when  $L = 1$ , but are omitted for brevity of exposition. Moreover, in that case, when Assumption 2.2 does not hold, we may again identify the above parameters for the population such that  $\Theta \in \text{supp}(f_{V|X}(\cdot|x))$  and  $X = x$ .

**Theorem 3.5.** Let Assumptions 2.1 (with (1)), 3.2, 3.8 and 3.9 hold. In case  $L = 1$ , for every  $\delta$  in  $\mathbb{R}$ , almost surely in  $x$  in  $\text{supp}(X)$ , for every  $\theta$  in  $\text{supp}(f_{\Theta|X}(\cdot|x))$ , we obtain

$$(3.15) \quad \text{UCDITE}(\delta, \theta, x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-it\delta} \frac{\overline{\partial_v \mathbb{E}[e^{itY} D | V = \cdot, X = x]}(\theta)}{\overline{\partial_v \mathbb{E}[e^{itY} (D - 1) | V = \cdot, X = x]}(\theta)} dt;$$

while, in case  $L \geq 2$ , for every  $\delta$  in  $\mathbb{R}$ , almost surely in  $x$  in  $\text{supp}(X)$ , for every  $\gamma$  in  $\text{supp}(f_{\tilde{\Gamma}|X}(\cdot|x))$ , we obtain:

$$(3.16) \quad \text{UCDITE}(\delta, \gamma, x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-it\delta} \frac{R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ e^{itY} D \mid (\tilde{S}, \tilde{V}) = \cdot, X = x \right]} \right] (\gamma)}{R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ e^{itY} (D - 1) \mid (\tilde{S}, \tilde{V}) = \cdot, X = x \right]} \right] (\gamma)} dt.$$

Note that there is a slight abuse of notations in (3.15) and (3.16) because the numerators and denominators are only defined almost surely in  $\theta$ , respectively  $\gamma$ . Implicitly, as in Heckman, Smith and Clements (1997), this result can be used to derive a test for the validity of Assumption 3.2: if for some  $x$  and  $\theta$  (respectively  $\gamma$ )  $\text{UCDITE}(\delta, \theta, x)$  (respectively  $\text{UCDITE}(\delta, \gamma, x)$ ) fails to be a density, it is an indication that Assumption 3.2 is incorrect.

#### 4. ESTIMATION OF STRUCTURAL PARAMETERS

In this section we focus on the case where  $L \geq 2$ , the case where  $L = 1$  requires minor modifications. The estimators that we consider in this section are not based on theoretical formulas that would only involve values at “infinity” of the instruments in the selection equation, nor situations with unselected samples. On the contrary, all the estimators make an efficient use of the data and use the whole sample. They involve trimming of tail values of the instruments to reduce the variance of our estimators and alleviate the potential influence of outliers on the estimation. We provide their detailed asymptotic properties in Section 5.

**4.1. Estimation of Quantities Related to Marginals.** We consider the following regularized inverse<sup>8</sup> of the Radon transform

$$(4.1) \quad A_T[f](\gamma) := \int_{H^+} \int_{-\infty}^{\infty} K_T(s'\gamma - u) f(s, u) du d\sigma(s)$$

where  $\sigma$  denotes the classical spherical measure<sup>9</sup> on the sphere  $\mathbb{S}^{L-1}$  and

$$(4.2) \quad \forall u \in \mathbb{R}, K_T(u) := 2(2\pi)^{-L} \int_0^{\infty} \cos(tu) t^{L-1} \psi\left(\frac{t}{T}\right) dt$$

<sup>8</sup>Up to our knowledge this modification of the classical Radon inverse has not been studied earlier in the literature.

<sup>9</sup>Its mass is the area of the sphere.

where  $T$  is a smoothing parameter and  $\psi$  is a symmetric rapidly decaying function in the Schwartz class

$$\mathcal{S}(\mathbb{R}) = \left\{ f \in C^\infty(\mathbb{R}) : \forall \alpha, \beta \in \mathbb{N}, |x|^\alpha \left| \partial^\beta f(x) \right| \xrightarrow{|x| \rightarrow \infty} 0 \right\},$$

such that  $\psi(0) = 1$ . For simplicity of exposition we will take  $\psi = \psi_0$  where  $\psi_0 : x \mapsto \exp\left(1 - \max\left\{\frac{1}{1-x^2}, 0\right\}\right)$  which has also support in  $[-1, 1]$ .

The previous identification sections suggest using this regularized inverse as a building block for a sample counterparts estimator of the inverse Radon transform of, say, a derivative of a regression function, *i.e.*,

$$A_{T_N} \left[ \widehat{\overline{\partial_v \mathbb{E}[\phi(Y)_\zeta(D) | (\tilde{S}, \tilde{V}) = \cdot, X = x]}} \right] (\gamma)$$

where  $\zeta(D)$  is either  $D$  or  $D - 1$  and  $\widehat{\overline{\partial_v \mathbb{E}[\phi(Y)_\zeta(D) | (\tilde{S}, \tilde{V}) = (\cdot, X = x]}}$  is the extension as 0 outside  $\text{supp}(\tilde{S}, \tilde{V})$  of an estimator of the derivative of the regression function (*e.g.*, using local polynomials). Moreover,  $T_N$  is chosen adequately, and tends to infinity as  $N$  goes to infinity. Replacing the various inverse Radon transforms by these regularized sample counterparts yields the following set of estimators:

$$(4.3) \quad \widehat{\overline{\text{UCATE}(\gamma, x) f_{\tilde{\Gamma}}(\gamma)}} = A_{T_N} \left[ \widehat{\overline{\partial_v \mathbb{E}[Y | (\tilde{S}, \tilde{V}) = \cdot, X = x]}} \right] (\gamma),$$

$$(4.4) \quad \mathcal{F}_1 \left[ \widehat{\overline{f_{Y_0 + \Delta, \tilde{\Gamma} | X}(\cdot | x)}} \right] (t, \gamma) = A_{T_N} \left[ \widehat{\overline{\partial_v \mathbb{E}[e^{itY} D | (\tilde{S}, \tilde{V}) = \cdot, X = x]}} \right] (\gamma),$$

$$(4.5) \quad \mathcal{F}_1 \left[ \widehat{\overline{f_{Y_0, \tilde{\Gamma} | X}(\cdot | x)}} \right] (t, \gamma) = A_{T_N} \left[ \widehat{\overline{\partial_v \mathbb{E}[e^{itY} (D - 1) | (\tilde{S}, \tilde{V}) = \cdot, X = x]}} \right] (\gamma),$$

$$(4.6) \quad \widehat{\overline{f_{\tilde{\Gamma} | X}(\gamma | x)}} = \max \left\{ A_{T_N} \left[ \widehat{\overline{\partial_v \mathbb{E}[D | (\tilde{S}, \tilde{V}) = \cdot, X = x]}} \right] (\gamma), 0 \right\}^{10}.$$

These individual elements can now be used to estimate many of the previously discussed quantities. From the first estimator, we may, for example, construct an estimator of the ATE.

$$\widehat{\text{ATE}}(x) = \int_{\mathbb{R}^L} \widehat{\overline{\text{UCATE}(\gamma, x) f_{\tilde{\Gamma} | X}(\gamma | x)}} \mathbf{1}\{\gamma \in \mathcal{B}_N\} d\gamma = \int_{\mathbb{R}^L} A_{T_N} \left[ \widehat{\overline{\partial_v \mathbb{E}[Y | (\tilde{S}, \tilde{V}) = \cdot, X = x]}} \right] (\gamma) \mathbf{1}\{\gamma \in \mathcal{B}_N\} d\gamma,$$

where  $\mathcal{B}_N$  is a large enough closed set in  $\mathbb{R}^L$  which diameter depends on the sample size.

Alternative estimators that circumvent the numerical integration in (4.1) can be obtained as follows. Start out by defining

$$\forall u \in \mathbb{R}, \tilde{K}_T(u) := -2(2\pi)^{-L} \int_0^T \sin(tu) t^L \psi\left(\frac{t}{T}\right) dt.$$

where  $K'_T(u) = \tilde{K}_T(u)$ , with  $K_T$  defined by (4.2) and

$$(4.7) \quad B_T[f](\gamma) := \int_{H^+} \int_{-\infty}^{\infty} \tilde{K}_T(s'\gamma - u) f(s, u) du d\sigma(s).$$

Because  $K_T$  is the Fourier transform of a function in  $\mathcal{S}(\mathbb{R})$ ,  $K_T$  and  $\tilde{K}_T$  belong to  $\mathcal{S}(\mathbb{R})$ . Thus, they decay to zero faster than any polynomial and belong to any  $L^p(\mathbb{R})$ <sup>11</sup>.

In addition, for the alternative estimators, we require the following assumption.

**Assumption 4.1.** (i) For almost every  $s \in H^+$  and almost surely in  $x \in \text{supp}(X)$ ,  
 $\text{supp}\left(f_{\tilde{V}|\tilde{S},X}(\cdot|s, x)\right) = \mathbb{R}$ .

(ii) For the function  $\phi$  considered, for almost every  $s \in H^+$  and almost surely in  $x \in \text{supp}(X)$ ,  $v \mapsto \mathbb{E}[\phi(Y)_\zeta(D)|(\tilde{S}, \tilde{V}) = (s, v), X = x]$  is continuous and  $v \mapsto \mathbb{E}[\phi(Y)_\zeta(D)|(\tilde{S}, \tilde{V}) = (s, v), X = x]$  and  $v \mapsto \partial_v \mathbb{E}[\phi(Y)_\zeta(D)|(\tilde{S}, \tilde{V}) = (s, v), X = x]$  are bounded by a polynomial in  $v$ .

This assumption allows an integration by parts argument for the regularized inverse, which produces a structure that is easier to implement. It is based on the following proposition.

**Proposition 4.1.** Under Assumption 4.1, for every  $T \in \mathbb{R}$  and  $\gamma \in \mathbb{R}^L$ ,

$$(4.8) \quad A_T \left[ \partial_v \mathbb{E} \left[ \phi(Y)_\zeta(D) | (\tilde{S}, \tilde{V}) = (\cdot), X = x \right] \right] (\gamma) = B_T \left[ \mathbb{E} \left[ \phi(Y)_\zeta(D) | (\tilde{S}, \tilde{V}) = (\cdot), X = x \right] \right] (\gamma)$$

$$(4.9) \quad = \mathbb{E} \left[ \frac{\tilde{K}_T(\tilde{S}'\gamma - \tilde{V}) \phi(Y)_\zeta(D)}{f_{\tilde{S}, \tilde{V}|X}(\tilde{S}, \tilde{V}|x)} \middle| X = x \right].$$

Equation (4.8) suggests that one can take as an estimator  $B_T$  applied to an estimator of the regression function. (4.9) suggests the following trimmed sample counterpart estimator

$$(4.10) \quad \frac{1}{N} \sum_{i=1}^N \frac{\tilde{K}_{T_N}(\tilde{s}'_i \gamma - \tilde{v}_i) T_{\tau_N}(\phi(y_i))_\zeta(d_i)}{\max\left(\widehat{f_{\tilde{S}, \tilde{V}|X}}(\tilde{s}_i, \tilde{v}_i, x), m_N\right)} \mathcal{K}_{\eta_N}(x_i - x)$$

where  $\widehat{f_{\tilde{S}, \tilde{V}|X}}$  is a plug-in estimator for  $f_{\tilde{S}, \tilde{V}|X}$ ,  $m_N$  is a trimming factor and  $\mathcal{K}_\eta$  is a standard multivariate kernel with bandwidth vector  $\eta_N$ ,  $T_\tau$  is defined for  $\tau$  positive and  $x$  in  $\mathbb{R}$  by

$$T_\tau(x) = -\tau \mathbf{1}\{x < -\tau\} + x \mathbf{1}\{|x| \leq \tau\} + \tau \mathbf{1}\{x > \tau\}$$

$T_N$ ,  $\tau_N$ ,  $m_N^{-1}$  and  $\eta_N^{-1}$  go to infinity as  $N$  goes to infinity. Trimming is introduced to avoid dividing by quantities that are too close to zero and thus giving a large weight to tail values of the instruments.

<sup>11</sup>This is a very nice feature that is not shared for the classical Radon inverse where  $\psi(x) = \mathbf{1}\{x \in [-1, 1]\}$ , also these yield good approximation results for the target quantities in every  $L^p(\mathbb{R}^L)$ .

Recall that this estimator can only be computed when  $f_{\tilde{S}, \tilde{V}|X}(\cdot|x)$  has full support, almost surely in  $x$  in  $\text{supp}(X)$ , implying that the density decays to zero both for large  $v$  and commonly when  $s$  approaches the boundary of  $H^+$ <sup>12</sup>. In particular, the true density  $f_{\tilde{S}, \tilde{V}, X}$  is not bounded away from zero. Moreover we are replacing  $f_{\tilde{S}, \tilde{V}, X}$  by an estimator, implying that the denominator could also be small due to estimation error. The truncation parameter  $\tau_N$  is useful when  $\phi$  is unbounded and  $\phi(Y_0)$  (if  $\varsigma(D) = D - 1$ ) or  $\phi(Y_1)$  (if  $\varsigma(D) = D$ ) have fat tails. In the absence of conditioning on covariates  $X$ , the estimator simplifies to

$$(4.11) \quad \frac{1}{N} \sum_{i=1}^N \frac{\tilde{K}_{T_N}(\tilde{s}'_i \gamma - \tilde{v}_i) T_{\tau_N}(\phi(y_i)) \varsigma(d_i)}{\max\left(\widehat{f_{\tilde{S}, \tilde{V}}(\tilde{s}_i, \tilde{v}_i)}, m_N\right)}.$$

Note the parallels, when  $\phi$  is the identity, between this estimator and the one in Lewbel (2007)<sup>13</sup> or Horvitz and Thompson (1952).

**4.2. Estimation of UCDITE and of the Distribution of Treatment Effects.** An estimator for UCDITE is obtained by applying the same principles. First, one may compute the following integral (4.12)

$$\widehat{\text{UCDITE}(\delta, \gamma, x)} = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-it\delta} K(th_{N,\gamma}) \frac{\mathcal{F}_1 \left[ \widehat{f_{Y_0+\Delta, \tilde{\Gamma}|X}(\cdot|x)} \right] (t, \gamma)}{\mathcal{F}_1 \left[ \widehat{f_{Y_0, \tilde{\Gamma}|X}(\cdot|x)} \right] (t, \gamma)} \mathbf{1} \left\{ \left| \mathcal{F}_1 \left[ \widehat{f_{Y_0, \tilde{\Gamma}|X}(\cdot|x)} \right] (t, \gamma) \right| > t_{N,t,\gamma} \right\} dt$$

where  $N$  is the sample size,  $h_{N,\gamma}$  is a bandwidth going to zero with  $N$ ,  $K$  denotes a kernel,  $t_{N,t,\gamma}$  a proper trimming factor.

A typical example of a kernel is  $K(t) = \mathbf{1}\{|t| \leq 1\}$ , with  $h_{N,\gamma} = 1/R_{N,\gamma}^\Delta$  it amounts to truncation of high frequencies. Devroye (1989) uses  $\max(1 - |t|, 0)$  for the estimation with  $L^1(\mathbb{R})$  loss. It is also possible to take  $K = \psi$  where  $\psi$  belongs to  $\mathcal{S}(\mathbb{R})$  and is such that  $\psi(0) = 1$  like in Section 4.1. In Section 5 we take  $K(t) = \psi_0(t)$  with support in  $[-1, 1]$  (recall that  $\psi_0(t) = \exp\left(1 - \max\left\{\frac{1}{1-t^2}, 0\right\}\right)$ ). For every  $\gamma$ , the quantity  $\left| \mathcal{F}_1 \left[ \widehat{f_{Y_0, \tilde{\Gamma}|X}(\cdot|x)} \right] (t, \gamma) \right|$  in the denominator of (4.12), decays to zero when  $t$  goes to infinity<sup>14</sup>. A smoothing kernel  $K$  should put less weight (possibly zero weight) on high frequencies<sup>15</sup> because this is where the denominator is small in modulus and the variance of the

<sup>12</sup>This is because  $\tilde{S}$  is a rescaled vector of original instruments and assuming that the density does not decay to zero when  $s$  approaches the boundary of  $H^+$  would be a very strong assumption on the tails of the distributions of the original instruments (see, *e.g.*, Beran, Feuerverger and Hall (1996) and Hoderlein, Klemelä and Mammen (2011)).

<sup>13</sup>Compared to the estimator (5) in Lewbel (2007), (4.11) has the extra weight  $\tilde{K}_{T_N}(\tilde{s}'_i \gamma - \tilde{v}_i)$  coming from the inversion of the Radon transform.

<sup>14</sup>This is a consequence of the Riemann-Lebesgue lemma.

<sup>15</sup>The parameter  $t$  is the frequency.



estimator can blow-up. In theory the bandwidth  $h_{N,\gamma}$  may depend on  $\gamma$ <sup>16</sup> because of the possible different rates of decay of  $\left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}|X}(\cdot|x) \right] (t, \gamma) \right|$  for different values of  $\gamma$ .

Compared to usual deconvolution estimators, (4.12) also has an indicator function with a trimming factor. It is used because we are considering a case where the denominator is unknown and has to be estimated. For that reason the denominator can also be small due to estimation error. This is the same structure as the estimator proposed in Neumann (1997)<sup>17</sup>.

In the simulation study in Section 5, we take  $h_{N,\gamma}$  and  $t_{N,t,\gamma}$  that do not depend on  $t$  and  $\gamma$  for the smoothing parameters in the estimators of the numerator and denominator in (4.12). Using  $K = \psi_0$  we also obtain exactly the same graphs for (4.12) regardless of whether or not we trim, so we present results with  $t_{N,t,\gamma} = 0$ . This is an important feature in practice, because in Proposition 4.5 we impose to tune  $t_{N,t,\gamma} = r_{Y_0,N}$ . But  $r_{Y_0,N}$  is unknown since it depends on the smoothness of  $f_{Y_0, \tilde{\Gamma}|X}$  which is impossible to estimate. We believe that this is just a technical issue and that in practice no trimming works perfectly fine at least for smooth  $f_{Y_0, \tilde{\Gamma}|X}$  (super smooth in our simulation example because it is Gaussian).

An estimator of the unconditional distribution of treatment effects uses as plug-ins (4.12) and an estimator of the mixing density  $f_{\tilde{\Gamma}|X}$ ; it is given by

$$(4.13) \quad \widehat{f_{\Delta|X}}(\delta|x) = \int_{\mathbb{R}^L} \widehat{\text{UCDITE}}(\delta, \gamma, x) \widehat{f_{\tilde{\Gamma}|X}}(\gamma|x) \mathbf{1}\{\gamma \in \mathcal{B}_N\} d\gamma.$$

**4.3. Rates of Convergence.** To analyze the rate of convergence of our estimators, we have to introduce additional notation. We denote by  $\|\cdot\|_p$  for  $p \in [1, \infty)$  the classical  $L^p$  norms and by  $\|\cdot\|_\infty$  the essential supremum norm, also called sup-norm. Moreover, for ease of notation we again omit the conditioning on  $X$  in this section.

**4.3.1. Estimation of Generic Quantities.** In this section we consider the estimation of one of the plug-in terms of the form

$$g(\gamma) = R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ \phi(Y) \varsigma(D) \mid (\tilde{S}, \tilde{V}) = \cdot \right]} \right] (\gamma),$$

for some function  $\phi$  and  $\varsigma(D)$  is either  $D$  or  $1 - D$ , by an estimator of the form (4.11)

$$\hat{g}(\gamma) = \frac{1}{N} \sum_{i=1}^N \frac{\tilde{K}_{T_N}(\tilde{s}'_i \gamma - \tilde{v}_i) T_{\tau_N}(\phi(y_i)) \varsigma(d_i)}{\max \left( \widehat{f_{\tilde{S}, \tilde{V}}}(\tilde{s}_i, \tilde{v}_i), m_N \right)}$$

<sup>16</sup>Every smoothing parameter can also depend on  $x$  but we omit it ease of notations.

<sup>17</sup>In Neumann (1997) and Comte and Lacour (2011) because the characteristic function of  $Y_0$  is estimated at rate  $1/\sqrt{N}$ ,  $t_{N,t,\gamma}$  could be taken equal to  $N^{-1/2}$ , independent of  $t$  and  $\gamma$ .

where  $T_N$ ,  $\tau_N$  and  $m_N^{-1}$  increase with the sample size.

We restrict our attention to the sup-norm loss because this is what is later required in Assumption 4.3 and Proposition 4.5 for plug-in estimation. We specify the result to the case of ordinary smooth functions (Sobolev classes). We will work with the following Sobolev spaces of all locally integrable functions that have weak derivatives up to order  $s$  in  $\mathbb{N} \setminus \{0\}$  (see, *e.g.*, Evans (1998))

$$W^{s,\infty}(\mathbb{R}^L) := \{f \in L^\infty(\mathbb{R}^L) : \forall |\alpha| \leq s, \partial^\alpha f \in L^\infty(\mathbb{R}^L)\}$$

where  $\alpha \in \mathbb{N}^L$ ,  $|\alpha| := \sum_{l=1}^L \alpha_l$  and  $\partial^\alpha f := \prod_{l=1}^L \partial_l^{\alpha_l} f$  is the  $\alpha^{th}$ -weak partial derivative of  $f$ . It is equipped with the norm

$$\|f\|_{s,\infty} := \sum_{\alpha: |\alpha| \leq s} \|\partial^\alpha f\|_\infty.$$

We will consider the following Sobolev ellipsoids defined for  $M$  positive by

$$W^{s,\infty}(M) := \{f \in W^{s,\infty}(\mathbb{R}^L) : \|f\|_{s,\infty} \leq M\}.$$

In what follows,  $\mathcal{B}_N$  is a closed set in  $\mathbb{R}^L$  and we denote by  $d(\mathcal{B}_N)$  its diameter for the Euclidian norm.

**Proposition 4.2.** Let Assumption 4.1 hold, and assume moreover that

- (i)  $g$  belongs to  $W_\infty^s(M)$  for some  $s$  in  $\mathbb{N} \setminus \{0\}$  and  $M$  positive ;
- (ii) there exists  $\alpha$  positive such that  $\log(T_N^3/m_N) + L \log(d(\mathcal{B}_N)) \leq \alpha$  ;
- (iii) there exists a sequence  $r_{IV,N}$  going to 0 as  $N$  goes to infinity and  $M_{IV}$  positive such that with probability one

$$(4.14) \quad \overline{\lim}_{N \rightarrow \infty} r_{IV,N}^{-1} \max_{i=1,\dots,N} \left| f_{\tilde{S}, \tilde{V}}(\tilde{s}_i, \tilde{v}_i) - \widehat{f_{\tilde{S}, \tilde{V}}}(\tilde{s}_i, \tilde{v}_i) \right| \leq M_{IV} ;$$

then for some constants  $M(\alpha)$  (which only depends on  $\alpha$  and  $L$ ) and  $C(s)$  (which only depends on  $s$  and  $\psi$ ), with probability one, for every  $\epsilon$  positive, for  $N$  large enough

$$\begin{aligned} \|(\hat{g} - g) \mathbf{1}\{\mathcal{B}_N\}\|_\infty &\leq (M_{IV} + \epsilon) \min(\tau_N, \|\phi\|_\infty) r_{IV,N} m_N^{-1} \left\| \mathbb{E} \left[ \frac{|\tilde{K}_{T_N}(\tilde{S}'\gamma - \tilde{V})|}{\max(\widehat{f_{\tilde{S}, \tilde{V}}}(\tilde{S}, \tilde{V}), m_N)} \right] \right\|_\infty \\ &\quad + (M_{IV} + \epsilon) \min(\tau_N, \|\phi\|_\infty) r_{IV,N} m_N^{-3/2} (M(\alpha) + \epsilon) \left( \frac{\log N}{N} \right)^{1/2} T_N^{L+1/2} \\ &\quad + m_N^{-1/2} (M(\alpha) + \epsilon) \min(\tau_N, \|\phi\|_\infty) \left( \frac{\log N}{N} \right)^{1/2} T_N^{L+1/2} \\ &\quad + \min(\tau_N, \|\phi\|_\infty) \sup_{\gamma \in \mathcal{B}_N} \int_{\{(s,v): f_{\tilde{S}, \tilde{V}}(s,v) < m_N\}} \left| \tilde{K}_{T_N}(s'\gamma - v) \right| d\sigma(s) dv \\ &\quad + MC(s) T_N^{-s} \end{aligned}$$

$$+ \frac{1}{(2\pi)^L} T_N^{L+2} \|t\|^L \psi \|_1 \mathbb{E}[|\phi(Y_j)| \mathbf{1}\{|\phi(Y_j)| > \tau_N\}]$$

where  $j = 1$  if  $\zeta(D) = D$  and  $j = 0$  if  $\zeta(D) = D - 1$ .

Let us make a few comments on this result.

- (1) When  $\text{supp}(\tilde{\Gamma})$  is bounded then we can take  $\mathcal{B}_N = \text{supp}(\tilde{\Gamma}) = \text{supp}(g)$ .
- (2) Condition (4.15) can be relaxed to "bounded in probability" if we simply want to prove convergence in probability. This is actually the only thing that we need for the properties of the estimation of UCDITE that will follow.
- (3) There are various ways to bound from above

$$\mathbb{E} \left[ \frac{|\tilde{K}_{T_N}(\tilde{S}'\gamma - \tilde{V})|}{\max(\widehat{f_{\tilde{S}, \tilde{V}}}(\tilde{S}, \tilde{V}), m_N)} \right].$$

A first uniform upper bound uses

$$\mathbb{E} \left[ \frac{|\tilde{K}_{T_N}(\tilde{S}'\gamma - \tilde{V})|}{\max(\widehat{f_{\tilde{S}, \tilde{V}}}(\tilde{S}, \tilde{V}), m_N)} \right] \leq \frac{|\mathbb{S}^{L-1}|}{2} \|\tilde{K}_{T_N}\|_1$$

where  $|\mathbb{S}^{L-1}|$  is the surface of the sphere  $\mathbb{S}^{L-1}$ . A second uniform upper bound, that has an analytic form<sup>18</sup> is given by

$$\mathbb{E} \left[ \frac{|\tilde{K}_{T_N}(\tilde{S}'\gamma - \tilde{V})|}{\max(\widehat{f_{\tilde{S}, \tilde{V}}}(\tilde{S}, \tilde{V}), m_N)} \right] \leq \mathbb{E} \left[ \left( \frac{|\tilde{K}_{T_N}(\tilde{S}'\gamma - \tilde{V})|}{\max(\widehat{f_{\tilde{S}, \tilde{V}}}(\tilde{S}, \tilde{V}), m_N)} \right)^2 \right]^{1/2} \leq m_N^{-1} T_N^{2L+1}.$$

The last inequality uses (7.5) in the appendix.

- (4) Because  $\tilde{K}_{T_N}$  is integrable<sup>19</sup>,  $\int_{\{(s,v): f_{\tilde{S}, \tilde{V}}(s,v) < m_N\}} |\tilde{K}_{T_N}(s'\gamma - v)| d\sigma(s) dv$  goes to zero as  $m_N$  goes to zero and the rate of convergence to zero depends on the tails of  $f_{\tilde{S}, \tilde{V}}$ . Fat tails imply that this term is small. It could be made equal to zero if  $f_{\tilde{S}, \tilde{V}}$  were bounded away from zero.
- (5) The contribution  $MC(s)T_N^{-s}$  is an upper bound on the approximation error for functions in the ellipsoid  $W^{s,\infty}(\mathbb{R}^L)$ .
- (6) We use truncation, because we deal with the fluctuation terms using the basic Bernstein inequality. It is possible to use the Bernstein inequality for random variables with bounded Orlicz norms (see, *e.g.*, Lemma 2.2.11 of Van der Vaart and Wellner (1996)), or other concentration

<sup>18</sup>We expect that it is not as sharp as the previous upper bound where, unfortunately, we do not have an upper bound for  $\|\tilde{K}_{T_N}\|_1$  in terms of its dependence in  $T_N$ .

<sup>19</sup>It belongs to  $\mathcal{S}(\mathbb{R})$ .

results and avoid truncation in certain cases. When  $\phi$  is bounded, we can take  $\tau_N = \|\phi\|_\infty$  and the last term in the above upper bound disappears. In Proposition 4.2, we considered the property of estimators involving truncation for the estimation of the numerator of UCATE and for the estimation of the conditional variance of treatment effect when  $|Y_0|$  and  $|Y_1|$  can take arbitrarily large values. Because  $\mathbb{E}[|\phi(Y_0)| + |\phi(Y_1)|] < \infty$ ,  $\mathbb{E}[|\phi(Y_1)|\mathbf{1}\{|\phi(Y_1)| > \tau_N\}]$  and  $\mathbb{E}[|\phi(Y_0)|\mathbf{1}\{|\phi(Y_0)| > \tau_N\}]$  go to zero when  $\tau_N$  goes to infinity.

In the ideal case where: (1)  $f_{\tilde{S}, \tilde{V}}$  is bounded away from zero, (2) its density is smooth enough for the first term to be negligible and (3) the bias due to truncation is negligible (*e.g.* when  $\phi$  is bounded), we obtain, for some  $M_I$ , with probability 1,

$$\overline{\lim}_{N \rightarrow \infty} \left( \frac{\log N}{N} \right)^{-\frac{s}{2s+2L+1}} \|\hat{g} - g\|_\infty \leq M_I$$

by taking  $T_N$  of the order of  $(N/\log(N))^{1/(2s+2L+1)}$ . Recall that the rate of direct density estimation would be  $(N/\log(N))^{s/(2s+L)}$ . This is important, because it says that the degree of ill-posedness due to the presence of the unbounded operator is  $\frac{L+1}{2}$ . Recall that in positron emission tomography, the classical statistical inverse problem involving the Radon transform, the degree of ill-posedness is  $\frac{L-1}{2}$  (see, *e.g.*, Korostelev and Tsybakov (1993)). Here we pay an additional price of 1 due to the extra differentiation. However, this degree of ill-posedness does not properly account for the difficulty of the problem. Equation (3.1) states that a regression function  $r$  is of the form  $r = Qf$  where  $Q$  is an operator which has an unbounded inverse. The quantity  $\frac{L-1}{2}$  only accounts for the smoothing properties of the operator  $K$ . But for identification we assumed that, in the original scale, all regressors but possibly  $V$  have full support. This implies that in most cases  $f_{\tilde{S}}(s)$  is not bounded away from zero and the rate of estimation of the regression function with  $L^\infty$  loss is slower than when the regressors have support on a compact set and their density is bounded from below. Estimation of a regression function when the density of the regressor can be 0 on its support (degeneracy) has been studied by several authors (see, *e.g.*, Hall, Marron, Neumann, Tetterington (1997), Guerre (1999) and Gaiffas (2005) and (2009))<sup>20</sup>. In our inverse problem setup this translates in the fact that in many cases  $\int_{\{(s,v): f_{\tilde{S}, \tilde{V}}(s,v) < m_N\}} \left| \tilde{K}_{T_N}(s'\gamma - v) \right| d\sigma(s)dv$  cannot be made negligible which implies that a second degree of ill-posedness also has to be taken into account. Finally, a third degree of ill-posedness can appear when the conditional distributions of  $\phi(Y_0)$  and/or  $\phi(Y_1)$  given  $\tilde{\Gamma}$  have heavy tails.

<sup>20</sup> Upper bounds in an inverse problem setting for specific estimators are given in Hoderlein, Klemelä and Mammen (2011) and Gautier and Kitamura (2009).

An analogue of Proposition 4.2 has already been established when  $\phi = 1$  in Gautier and Kitamura (2009) with the scaling of Section 6.1 and an estimator based on smoothed projection kernels in the Fourier domain. In this paper, we use a function  $\psi$  in  $\mathcal{S}(\mathbb{R})$  for the same reason for which we used smoothed projection kernels in Gautier and Kitamura (2009): to obtain rates of convergence for all  $L^p$  risks for  $1 \leq p \leq \infty$ <sup>21</sup>. In Gautier and Le Pennec (2011) this smoothed projection kernel is used together with the Littlewood-Paley decomposition and a quadrature formula to obtain a needlet estimator. Gautier and Le Pennec (2011) provide minimax lower bounds for the estimation when  $\phi = 1$  and show that their data-driven estimator is adaptive.

We have only considered the estimation of smooth functions for simplicity. If we consider “super smooth” functions, we expect that for certain functions  $\psi$  we could replace  $MC(s)T_N^{-s}$  by an exponentially small term. In that case, like in statistical deconvolution (see, *e.g.*, Butucea (2004) and Butucea and Tsybakov (2007)), for a nicely behaved density of the instruments and for  $\phi$  bounded, we could obtain parametric rates of convergence up to a logarithmic factor. Cavalier (2000) considers the estimation at a point of super smooth functions for the positron emission tomography problem in  $\mathbb{R}^L$ . The setup we consider is more involved because the inverse problem involves an extra derivative and we are in a regression framework with (1) random regressors whose density could be arbitrarily close to zero on its support and (2) possibly fat tails of the variables of interest. We do however consider super smooth functions in the more classical deconvolution framework in Section 4.3.4.

4.3.2. *Estimation of the Plug-in Terms for  $f_{\Delta|\tilde{\Gamma}}$ .* In this section we consider the estimation of the partial Fourier transforms which are used as plug-ins in Section 4.3.4. Denote by

$$g(\gamma, t) = R^{-1} \left[ \partial_v \mathbb{E} \left[ e^{itY} \zeta(D) \middle| \left( \tilde{S}, \tilde{V} \right) = \cdot \right] \right] (\gamma)$$

and consider an estimator of the form (4.11)

$$\hat{g}(\gamma, t) = \frac{1}{N} \sum_{i=1}^N \frac{\tilde{K}_{T_N}(\tilde{s}_i' \gamma - \tilde{v}_i) e^{ity_i} \zeta(d_i)}{\max \left( \widehat{f_{\tilde{S}, \tilde{V}}}(\tilde{s}_i, \tilde{v}_i), m_N \right)}$$

For technical reasons we introduce a maximum value  $R_N^{\max}$  for the inverse of the smoothing parameter  $h_{N,\gamma}^{-1}$ <sup>22</sup> in the estimator (4.12). In section 4.3.4 we will only be able to adjust  $h_{N,\gamma}^{-1}$  in the range  $[0, R_N^{\max}]$ .

**Proposition 4.3.** Make Assumption 4.1 and assume

<sup>21</sup>This is important to handle those plug-in terms to allow to use the whole range of the Hölder and Young inequalities.

<sup>22</sup>This is a classical feature of wavelet thresholding estimators and is called a maximal resolution level.

- (i) there exists  $s$  in  $\mathbb{N} \setminus \{0\}$  and  $M$  positive such that for every  $t$  in  $\mathbb{R}$ ,  $\Re[g(\gamma, t)]$  and  $\Im[g(\gamma, t)]$  belong to  $W_\infty^s(M)$  ;
- (ii) there exists  $\alpha$  positive such that  $\log(T_N^3/m_N) + \log(R_N^{\max}) + L \log(d(\mathcal{B}_N)) \leq \alpha$  ;
- (iii) there exists a sequence  $r_{IV,N}$  going to 0 as  $N$  goes to infinity and  $M_{IV}$  positive such that with probability one

$$(4.15) \quad \overline{\lim}_{N \rightarrow \infty} r_{IV,N}^{-1} \max_{i=1, \dots, N} \left| f_{\tilde{S}, \tilde{V}}(\tilde{s}_i, \tilde{v}_i) - \widehat{f_{\tilde{S}, \tilde{V}}}(\tilde{s}_i, \tilde{v}_i) \right| \leq M_{IV} ;$$

then for some constants  $M(\alpha)$  (which only depends on  $\alpha$  and  $L$ ) and  $C(s)$  (which only depends on  $s$  and  $\psi$ ), with probability one, for every  $\epsilon$  positive, for  $N$  large enough

$$\begin{aligned} & \sup_{t \in [-R_N^{\max}, R_N^{\max}], \gamma \in \mathcal{B}_N} |(\hat{g} - g)(t, \gamma)| \\ & \leq 2(M_{IV} + \epsilon) r_{IV,N} m_N^{-1} \left\{ \left\| \mathbb{E} \left[ \frac{|\tilde{K}_{T_N}(\tilde{S}'\gamma - \tilde{V})|}{\max(f_{\tilde{S}, \tilde{V}}(\tilde{S}, \tilde{V}), m_N)} \right] \right\|_\infty + m_N^{-1/2} (M(\alpha) + \epsilon) \left( \frac{\log N}{N} \right)^{1/2} T_N^{L+1/2} \right\} \\ & + 2m_N^{-1/2} (M(\alpha) + \epsilon) \left( \frac{\log N}{N} \right)^{1/2} T_N^{L+1/2} \\ & + 2 \sup_{\gamma \in \mathcal{B}_N} \int_{\{(s,v): f_{\tilde{S}, \tilde{V}}(s,v) < m_N\}} |\tilde{K}_{T_N}(s'\gamma - v)| d\sigma(s) dv \\ & + 2MC(s) T_N^{-s} \end{aligned}$$

This is the same upper bound as in Proposition 4.2 up to a factor 2 (we separate the real and imaginary part of  $e^{ity}$ ) and to a larger constant  $M(\alpha)$ , and the same remarks apply.

4.3.3. *Estimation of  $f_\Delta$ .* We start with a proposition that relates the estimation of  $f_\Delta$  to the estimation of  $f_{\tilde{\Gamma}}$  and of  $f_{\Delta|\tilde{\Gamma}}$ . To this end, we make the following assumptions.

**Assumption 4.2.**  $f_{\tilde{\Gamma}} \in L^\infty(\mathbb{R}^L)$  and there exists  $M_\Delta$  positive such that  $\sup_{\delta \in \mathbb{R}, \gamma \in \text{supp}(\tilde{\Gamma})} f_{\Delta|\tilde{\Gamma}}(\delta|\gamma) \leq M_\Delta$ .

In the next proposition we give an upper bound on the error in estimating  $f_\Delta$  when we use the estimator (4.13).

**Proposition 4.4.** Let Assumptions (1) and 4.2 hold, then for every measurable set  $\mathcal{B}_N$  in  $\mathbb{R}^L$ ,

$$\left\| \widehat{f_\Delta} - f_\Delta \right\|_2^2 \leq 3M_\Delta \left( \int_{\mathcal{B}_N^c} f_{\tilde{\Gamma}} d\gamma + M_\Delta \left\| \frac{\widehat{f_{\tilde{\Gamma}}} - f_{\tilde{\Gamma}}}{f_{\tilde{\Gamma}}} \mathbf{1}_{\{\mathcal{B}_N\}} \right\|_\infty^2 \right)$$

$$(4.16) \quad + 3 \|f_{\tilde{\Gamma}}\|_{\infty} \left( 1 + \left\| \frac{\widehat{f_{\tilde{\Gamma}}} - f_{\tilde{\Gamma}}}{f_{\tilde{\Gamma}}} \mathbf{1}_{\{\mathcal{B}_N\}} \right\|_{\infty} \right)^2 \left\| \left( \widehat{f_{\Delta|\tilde{\Gamma}}}(\cdot|\star) - \overline{f_{\Delta|\tilde{\Gamma}}}(\cdot|\star) \right) \mathbf{1}_{\{\star \in \mathcal{B}_N\}} \right\|_2^2.$$

When  $\text{supp}(\tilde{\Gamma})$  is bounded, and  $f_{\tilde{\Gamma}}$  is bounded away from zero on its support, then we can take  $\mathcal{B}_N = \text{supp}(\tilde{\Gamma})$ . Otherwise  $\mathcal{B}_N$  should be (1) small enough so that  $f_{\tilde{\Gamma}}$  is bounded away from zero on  $\mathcal{B}_N$  (recall as well that its diameter should not grow faster than polynomially in  $N$  to be able to apply the result from Section 4.3.1), and (2) large enough so that  $\mathbb{P}(\tilde{\Gamma} \in \mathcal{B}_N^c) = \int_{\mathcal{B}_N^c} f_{\tilde{\Gamma}} d\gamma$  is small.

The next section studies the convergence to zero of the term

$$\left\| \left( \widehat{f_{\Delta|\tilde{\Gamma}}}(\cdot|\star) - \overline{f_{\Delta|\tilde{\Gamma}}}(\cdot|\star) \right) \mathbf{1}_{\{\star \in \mathcal{B}_N\}} \right\|_2.$$

4.3.4. *Estimation of  $f_{\Delta|\tilde{\Gamma}}$ .* In order to work with smoothing and trimming factors in (4.12) that are independent of  $t$  and  $\gamma$ , we work with sup-norm consistency of the estimators of the partial Fourier transforms.

**Assumption 4.3.**

$$\sup_{t \in [-R_N^{\max}, R_N^{\max}], \gamma \in \mathcal{B}_N} \left| \mathcal{F}_1 \left[ \widehat{f_{Y_0+\Delta, \tilde{\Gamma}}} \right] - \mathcal{F}_1 \left[ f_{Y_0+\Delta, \tilde{\Gamma}} \right] \right| = O_p(r_{Y_0+\Delta, N})$$

$$\sup_{t \in [-R_N^{\max}, R_N^{\max}], \gamma \in \mathcal{B}_N} \left| \mathcal{F}_1 \left[ \widehat{f_{Y_0, \tilde{\Gamma}}} \right] - \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] \right| = O_p(r_{Y_0, N}).$$

Recall that  $R_N^{\max}$  is a maximal resolution level and we assume that  $h_{N, \gamma}^{-1} \leq R_N^{\max}$  and that  $\mathcal{B}_N$  is a domain in  $\mathbb{R}^L$  that could grow as  $N$  goes to infinity if  $\text{supp}(\tilde{\Gamma})$  is unbounded. Rates of estimation  $r_{Y_0+\Delta, N}$  and  $r_{Y_0, N}$  are given in Section 4.3.2.

Unlike deconvolution with noise observed on a preliminary sample, in this setup each rate is nonparametric; it is the rate of estimation in the respective inverse problems. Rates in sup norm are given in Proposition 4.2.

**Proposition 4.5.** Let Assumptions 2.1 (with (1)), 2.2, 3.2, 3.8, 3.9, 4.2 hold. Assume that  $K$  has support in  $[-1, 1]$  and that  $h_{N, \gamma}^{-1} \leq R_N^{\max}$ . Take  $t_{N, t, \gamma} = r_{Y_0, N}$ . The following upper bound holds

$$(4.17) \quad \left\| \left( \widehat{f_{\Delta|\tilde{\Gamma}}}(\cdot|\star) - \overline{f_{\Delta|\tilde{\Gamma}}}(\cdot|\star) \right) \mathbf{1}_{\{\star \in \mathcal{B}_N\}} \right\|_2^2$$

$$(4.18) \quad = O_p \left( \int_{\mathcal{B}_N \cap \text{supp}(\tilde{\Gamma})} \int_{-\infty}^{\infty} \left[ (1 - K(t h_{N, \gamma}))^2 \left| \mathcal{F}_1 \left[ f_{\Delta|\tilde{\Gamma}} \right] (t|\gamma) \right|^2 \right. \right.$$

$$\left. \left. + \frac{K(t h_{N, \gamma})^2}{\left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right|^2} \left( r_{Y_0+\Delta, N}^2 + \left| \mathcal{F}_1 \left[ f_{\Delta, \tilde{\Gamma}} \right] (t, \gamma) \right|^2 r_{Y_0, N}^2 \right) \right] dt d\gamma \right).$$

The first term in the upper bound is the square of the approximation bias. Consider now the following classes of ellipsoids for  $f_{\Delta|\tilde{\Gamma}}$

$$\mathcal{A}_{\delta,r,a}(L) = \left\{ f \in L^2(\mathbb{R}) : \int_{-\infty}^{\infty} |\mathcal{F}[f](t)|^2 (1+t^2)^\delta \exp(2a|t|^r) dt d\gamma \leq L^2 \right\}$$

where  $r \geq 0$ ,  $a > 0$ ,  $\delta \in \mathbb{R}$  and  $\delta > 1/2$  if  $r = 0$ ,  $l > 0$ . The case  $r > 0$  corresponds to an extension of the case of super smooth functions, otherwise the functions are extensions of ordinary smooth functions (in the Sobolev class). When  $K(t) = \mathbf{1}\{|t| \leq 1\}$ ,  $h_{N,\gamma}$  is of the form  $1/R_N^\Delta$ , and  $f$  belongs to  $\mathcal{A}_{\delta,r,a}(L)$  then we have

$$\int_{-\infty}^{\infty} (1 - K(t h_{N,\gamma}))^2 |\mathcal{F}[f](t)|^2 dt \leq L^2 \left( (R_N^\Delta)^2 + 1 \right)^{-\delta} \exp\left(-2a (R_N^\Delta)^r\right).$$

The next proposition considers the case when  $K(t) = \mathbf{1}\{|t| \leq 1\}$  and  $h_{N,\gamma}$  is of the form  $1/R_N^\Delta$ . We make the following assumption on the decay rate of  $\left| \mathcal{F}_1 \left[ f_{Y_0|\tilde{\Gamma}} \right] (t|\gamma) \right|$  that strengthens Assumption 3.9.

**Assumption 4.4.** There exists  $s \geq 0$ ,  $b > 0$ ,  $\eta \in \mathbb{R}$  ( $\eta > 0$  if  $s = 0$ ) and  $k_0, k_1 > 0$  such that for every  $\gamma$  in  $\text{supp}(\tilde{\Gamma})$ ,

$$k_0(1+t^2)^{-\eta/2} \exp(-b|t|^s) \leq \left| \mathcal{F}_1 \left[ f_{Y_0|\tilde{\Gamma}} \right] (t|\gamma) \right| \leq k_1(1+t^2)^{-\eta/2} \exp(-b|t|^s)$$

In the proposition below we use the short hand notation  $\lambda_N = \lambda \left( \mathcal{B}_N \cap \text{supp}(\tilde{\Gamma}) \right)$ , where  $\lambda(B)$  is the Lebesgue measure of a set  $B$ .

**Proposition 4.6.** Let Assumptions 2.1 (with (1)), 2.2, 3.2, 3.8, 3.9, 4.2 and 4.3 and 4.4. Assume for every  $\gamma \in \text{supp}(\tilde{\Gamma})$ ,  $f_{\Delta|\tilde{\Gamma}}(\cdot|\gamma)$  belongs to  $\mathcal{A}_{\delta,r,a}(L)$ . The following upper bounds hold for every  $R_N^\Delta \leq R_N^{\max}$  and  $\mathcal{B}_N$  measurable set in  $\mathbb{R}^L$ ,

(C-1) if  $s = r = 0$ , then

$$\begin{aligned} & \lambda_N^{-1} \left\| \left( \widehat{f_{\Delta|\tilde{\Gamma}}}(\cdot|\star) - \overline{f_{\Delta|\tilde{\Gamma}}}(\cdot|\star) \right) \mathbf{1}\{\star \in \mathcal{B}_N\} \right\|_2^2 \\ &= O_p \left( (R_N^\Delta)^{-2\delta} + r_{Y_0+\Delta,N}^2 (R_N^\Delta)^{2\eta+1} + r_{Y_0,N}^2 (R_N^\Delta)^{2\max(\eta-\delta,0)+1} \right); \end{aligned}$$

(C-2) if  $s > 0$  and  $r = 0$ ,

$$\begin{aligned} & \lambda_N^{-1} \left\| \left( \widehat{f_{\Delta|\tilde{\Gamma}}}(\cdot|\star) - \overline{f_{\Delta|\tilde{\Gamma}}}(\cdot|\star) \right) \mathbf{1}\{\star \in \mathcal{B}_N\} \right\|_2^2 \\ &= O_p \left( (R_N^\Delta)^{-2\delta} + e^{2b(R_N^\Delta)^s} \left( r_{Y_0+\Delta,N}^2 (R_N^\Delta)^{2\eta+1-s} + r_{Y_0,N}^2 (R_N^\Delta)^{\min(1+2\eta-s, 2(\eta-\delta))} \right) \right); \end{aligned}$$



(C-3) if  $s = 0$  and  $r > 0$ , then

$$\lambda_N^{-1} \left\| \left( \widehat{f_{\Delta|\tilde{\Gamma}}}(\cdot|\star) - \overline{f_{\Delta|\tilde{\Gamma}}}(\cdot|\star) \right) \mathbf{1}_{\{\star \in \mathcal{B}_N\}} \right\|_2^2 = O_p \left( (R_N^\Delta)^{-2\delta} e^{-2a(R_N^\Delta)^r} + r_{Y_0+\Delta,N}^2 (R_N^\Delta)^{2\eta+1} + r_{Y_0,N}^2 \right);$$

(C-4) if  $s > 0$  and  $r > 0$ , then

$$\begin{aligned} & \lambda_N^{-1} \left\| \left( \widehat{f_{\Delta|\tilde{\Gamma}}}(\cdot|\star) - \overline{f_{\Delta|\tilde{\Gamma}}}(\cdot|\star) \right) \mathbf{1}_{\{\star \in \mathcal{B}_N\}} \right\|_2^2 \\ &= O_p \left( (R_N^\Delta)^{-2\delta} e^{-2a(R_N^\Delta)^r} + r_{Y_0+\Delta,N}^2 (R_N^\Delta)^{2\eta+1-s} e^{2b(R_N^\Delta)^s} + r_{Y_0,N}^2 \Delta (R_M^\Delta) \right) \end{aligned}$$

where

$$\begin{aligned} \Delta (R_M^\Delta) &= (R_N^\Delta)^{\min(1+2\eta-s, 2(\eta-\delta))} e^{2b(R_N^\Delta)^s} \mathbf{1}_{\{s > r\}} + (R_N^\Delta)^{\max(2(\eta-\delta), 0)} e^{2(b-a)(R_N^\Delta)^s} \mathbf{1}_{\{r = s, b \geq a\}} \\ &+ \mathbf{1}_{\{r > s\} \cup \{r = s, b < a\}}. \end{aligned}$$

When  $\text{supp}(\tilde{\Gamma})$  is bounded then  $\lambda_N$  is a constant.

We did not want to include results in other norms than the  $L^2$  norm for brevity of exposition. However it is possible to obtain rates of convergence in sup-norm for an estimator with  $K = \psi$ , where  $\psi$  belongs to  $\mathcal{S}(\mathbb{R})$  and  $\psi(0) = 1$ . It would also be possible to use the sup-norm adaptive<sup>23</sup> estimator of Lounici and Nickl (2011).

## 5. SIMULATION STUDY

We consider the model (1.1)-(1.2) with the added specification

$$Y_0 = 1 + 1.5\Gamma + \Theta + \varepsilon_0$$

$$Y_1 = 3 + 2.5\Gamma - \Theta + \varepsilon_1$$

where  $(\Gamma, \Theta, \varepsilon_0, \varepsilon_1) = (1, -0.5, 0, 0) + W$ ,  $W$  is a centered Gaussian random vector with covariance matrix

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix},$$

$\tilde{S} = (\cos N_t, \sin N_t)$  where  $N_t$  is a truncated Gaussian random variable with mean  $\pi/2$  and variance  $\pi^2/16$  on the interval  $[0, \pi]$ ,  $V$  is a Gaussian random variable with mean -0.2 and variance 4 and  $V, N_t$

<sup>23</sup> Specific to the case where the distribution of the error is known.

and  $(\Gamma, \Theta, \varepsilon_0, \varepsilon_1)$  are independent. The sample size considered in the simulation study is  $N = 10\,000$  and we have performed  $S = 100$  Monte Carlo repetitions.

We present the results with the estimators (4.3)-(4.6), using  $\psi = \psi_0$ . Table 2 shows that, in this Monte-Carlo study, they clearly outperform the easier to calculate estimator (4.11). All numerical integrations were carried out by quadrature methods<sup>24</sup>. The choice of the smoothing parameters for the estimation of  $f_{\Gamma, \Theta}$  were  $T_N = 6$  for the regularized Radon inverse and  $h_N = 1$  for the bandwidth of the local polynomial estimator of  $\partial_v \mathbb{E} \left[ D \mid (\tilde{S}, \tilde{V}) = \cdot \right]$ , while we took  $T_N = 10$  and  $h_N = 1$  for the estimation of  $\widehat{\text{UCATE}} \times f_{\Gamma, \Theta}$  and did not use truncation.

In Figures 3 and 4 we compare the truth, an empirical average of estimators over  $S = 100$  simulations, and one typical simulation.

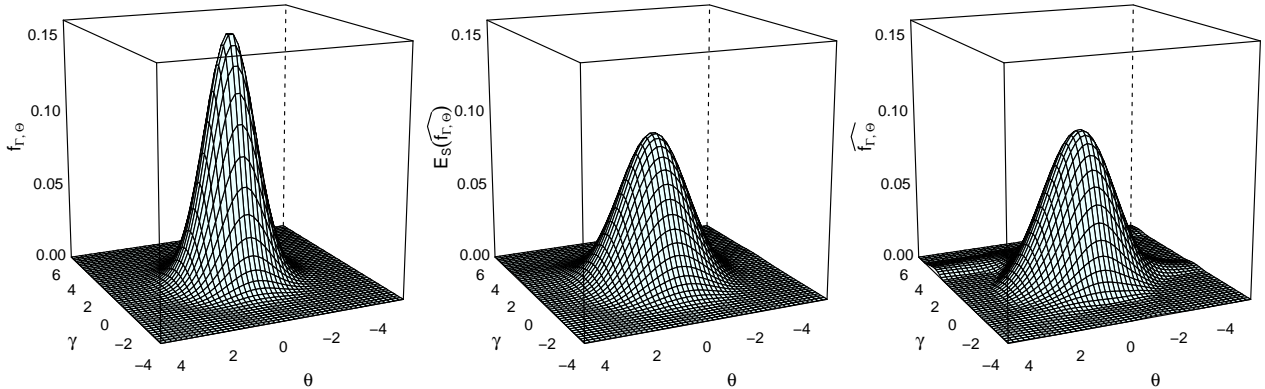


FIGURE 3. True  $f_{\Gamma, \Theta}$  (left), average over  $S$  replications of the estimator (middle) and the estimator calculated for one data set (right).

Figure 3, based on the rectangle  $[-4, 6] \times [-5.5, 4.5]$ , shows that most of the mass of  $\widehat{f_{\Gamma, \Theta}}$  is concentrated in a box around the mode at  $(1, -0.5)$ . Because UCATE is a conditional expectation given  $(\Gamma, \Theta)$ , it will only be well estimated at points where the density  $f_{\Gamma, \Theta}$  is not too low. Based on  $\widehat{f_{\Gamma, \Theta}}$  we experimented with several rectangular domains. Table 1 presents the empirical distribution

<sup>24</sup>We tried to apply an importance sampling Monte-Carlo method to calculate the multiple integrals using as proposals a Cauchy distribution for the integral with respect to  $v$  and a uniform distribution for the integral with respect to  $\phi$  ( $s = (\cos \phi, \sin \phi)$ ). Due to the presence of the  $\mathcal{S}(\mathbb{R})$  function  $\psi_0$  in  $K_T$  this Monte-Carlo approximation do have finite variance but the variance was too big to have a sufficiently good precision even using the highest possible sample size that we could generate in R. The variance of this Importance Sampling Monte-Carlo method is infinite if we replace  $\psi_0$  by an indicator function. It is possible that another choice of rapidly decreasing function  $\psi$  can make it feasible.

over  $S = 100$  replications of

$$\widehat{\text{ATE}}_j = \int_{\mathcal{B}_j} \text{UCATE} \times f_{\Gamma, \Theta}(\gamma, \theta) d\gamma d\theta$$

$$\widehat{\text{TT}}_j = \int_{\mathcal{B}_j} h_{\text{TT}}(\gamma) \text{UCATE} \times f_{\Gamma, \Theta}(\gamma, \theta) d\gamma d\theta$$

for estimators calculated on three such domains. The index  $j = 1$  corresponds to  $\mathcal{B}_1 = [-1.5, 3.5] \times [-3, 2]$  (2.5 standard errors in each direction), the index  $j = 2$  corresponds to  $\mathcal{B}_2 = [-1.75, 3.75] \times [-3.25, 2.25]$  (2.75 standard errors in each direction), while the index  $j = 3$  corresponds to  $\mathcal{B}_3 = [-2, 4] \times [-3.5, 2.5]$  (3 standard errors in each direction). For reference one should note that the true ATE is 4 while the true TT calculated via Monte-Carlo is 4.507 (0.0016).

	Mean	P5	P10	Median	P90	P95
$\widehat{\text{ATE}}_1$	3.91	2.98	3.28	3.88	4.67	4.81
$\widehat{\text{TT}}_1$	4.24	3.31	3.54	4.22	5.03	5.13
$\widehat{\text{ATE}}_2$	4.09	3.04	3.31	4.05	5.02	5.21
$\widehat{\text{TT}}_2$	4.46	3.37	3.61	4.42	5.36	5.61
$\widehat{\text{ATE}}_3$	4.22	2.89	3.21	4.22	5.42	5.69
$\widehat{\text{TT}}_3$	4.62	3.13	3.60	4.58	5.77	6.09

TABLE 1. The estimators indexed by 1 (resp. 2 and 3) correspond to the integration of  $\widehat{\text{UCATE}} \times f_{\Gamma, \Theta}$  on  $\mathcal{B}_1$  (resp.  $\mathcal{B}_2$  and  $\mathcal{B}_3$ ).

The plots in Figure 4 illustrate how our estimator performs for: the estimation of  $\text{UCATE} \times f_{\Gamma, \Theta}$  (top), and the estimation of UCATE only (bottom). For the former, we have used the domain  $\mathcal{B}_1$ , while for the latter we have employed the smaller domain  $\mathcal{B}_4 = [-0.5, 2.5] \times [-2, 1]$ . Note that  $\text{UCATE} \times f_{\Gamma, \Theta}$  should always be much more difficult to estimate than  $f_{\Gamma, \Theta}$  because it involves a regression function with a conditional expectation with respect to  $(\Gamma, \Theta)$  and the density of  $(\Gamma, \Theta)$  is not bounded away from zero. In the same spirit, UCATE should be even more difficult to estimate in this simulation setup as the tails of the numerator are fatter than that of the denominator. However, this simulation study shows how estimators of the denominator and numerator of UCATE that do not perform extremely well when the risk is defined in terms of the sup-norm (see, *e.g.*, the heights of the peaks), can perform reasonably well when estimating UCATE.

On the left panel of Figure 5 we compare the true  $\text{UCDITE}(\delta, (1, -0.5))$  to an estimator with  $K(t) = \exp\left(1 - \max\left\{\frac{1}{1-t^2}, 0\right\}\right)$  and  $R_{N, \gamma}^{\Delta} = 0.85$ . On the right panel we compare the true  $f_{\Delta}$  to

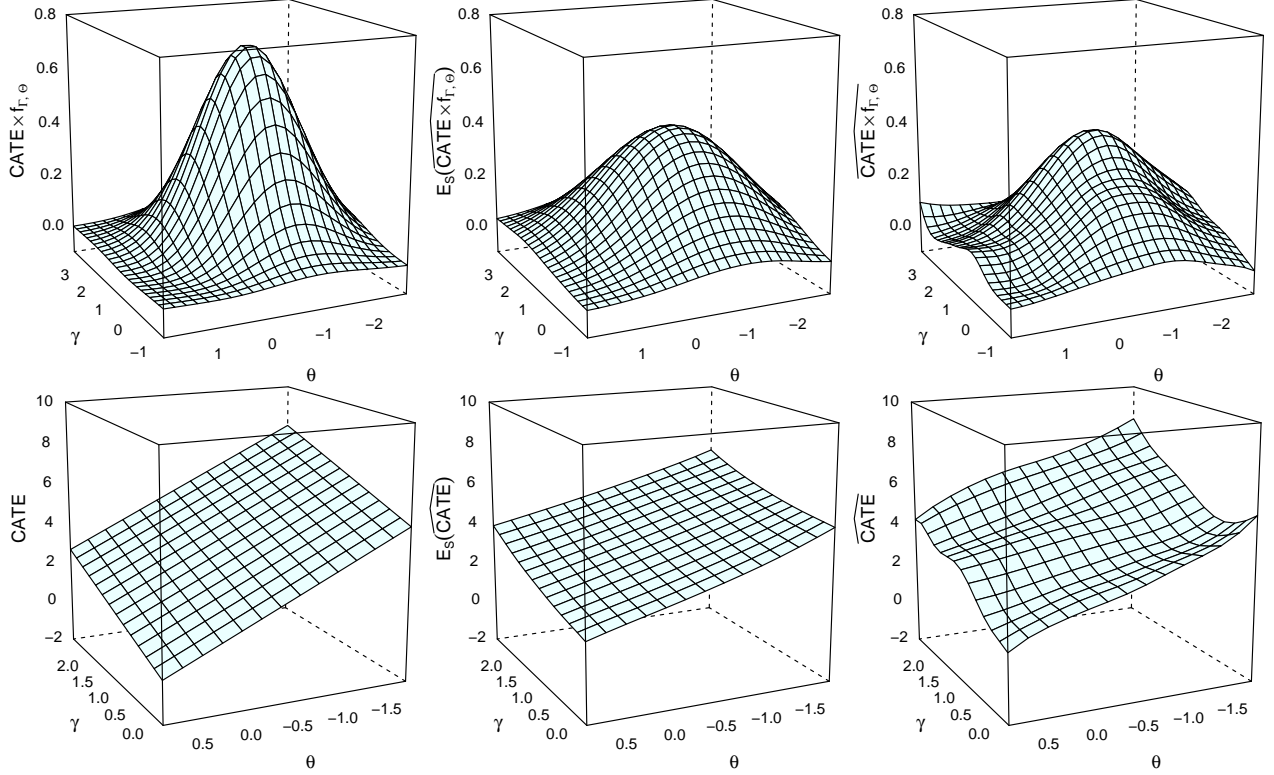


FIGURE 4.  $\text{UCATE} \times f_{\Gamma, \theta}$  (up) and  $\text{UCATE}$  (down), truth (left), average over  $S$  replications of the estimator (middle) and the estimator calculated for one data set (right).

Average Error	Estimators (4.3) and (4.6)	Estimator (4.11)
$g = f_{\Gamma, \theta}$		
$\frac{1}{S} \sum_{s=1}^S \ \hat{g}_s - g\ _{\infty}$	0.0635	0.0723
$\frac{1}{S} \sum_{s=1}^S \ \hat{g}_s - g\ _2$	0.00227	0.184
$g = \text{UCATE} \times f_{\Gamma, \theta}$		
$\frac{1}{S} \sum_{s=1}^S \ \hat{g}_s - g\ _{\infty}$	0.314	0.396
$\frac{1}{S} \sum_{s=1}^S \ \hat{g}_s - g\ _2$	0.000221	0.818

TABLE 2. Comparison between the estimators (4.6) and (4.3) and (4.11) for the numerator and denominator of  $\text{UCATE}$ . The smoothing parameter in (4.11) is  $T_N = 1.8$  for the estimation of  $f_{\Gamma, \theta}$  and  $T_N = 1.7$  for the estimation of  $\text{UCATE} \times f_{\Gamma, \theta}$ . It has been adjusted to perform as well as possible in sup-norm. The integration for the calculation of the  $L^2$ -norm is carried out on the domain  $\mathcal{B}_1$ .

an estimator obtained via a numerical integration on the box  $\mathcal{B}_1$ , with the same choice of  $R_{N,\gamma}^\Delta$  and without trimming. Indeed, we tried several values of a trimming parameter and all estimates were virtually indistinguishable. Both estimators are calculated on the sample that we present on the right panel in figures 3 and 4.

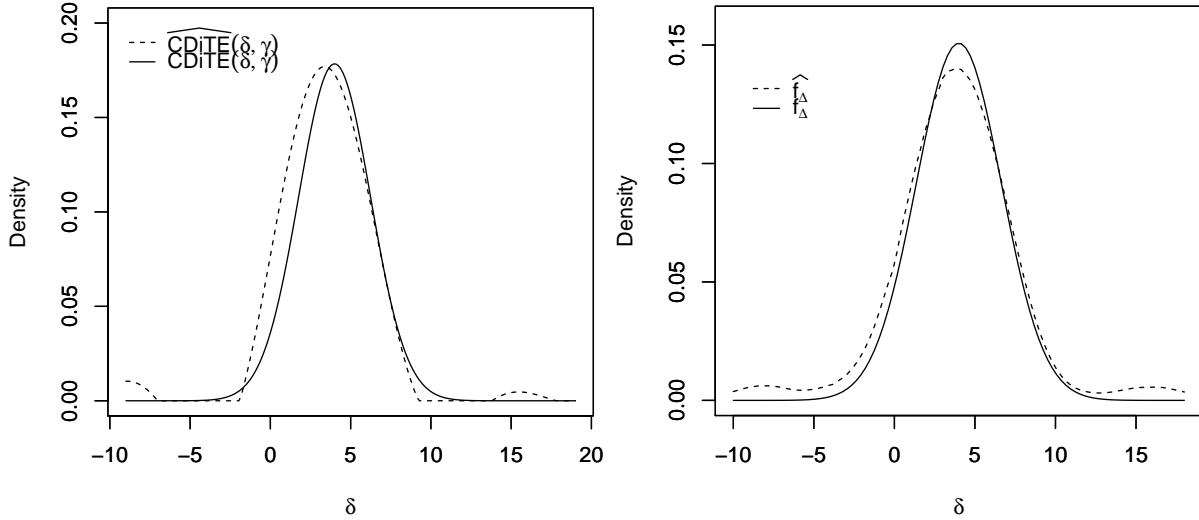


FIGURE 5. Comparison between (1) an estimator of UCDITE and the truth calculated at  $(1, -0.5)$ , the mode of  $(\Gamma, \Theta)$  (left) and (2) an estimator of  $f_\Delta$  and the truth (right).

In this simulation study we observe that the estimation of UCATE is paradoxically more difficult than the estimation of UCDITE. Indeed, in our setup the denominator of UCATE has thinner tails than the numerator. It is the opposite for UCDITE where the denominator has fatter tails than the numerator. Also, for UCDITE, the target density is super smooth and it is known that we can obtain very good rates of convergence for deconvolution estimators in that case (see the rates (A-4) in Proposition 4.6). Most quantities of interest are conditional expectations given  $(\Gamma, \Theta)$ . For these quantities it is vain to attempt to obtain an estimate at points where the density of  $(\Gamma, \Theta)$  is too low. We recommend to start by drawing plots of the density of  $(\Gamma, \Theta)$ , and ruling out such areas<sup>25</sup>. It is also more trustworthy to plot UCATE on a small domain. A graphical representation of UCATE allows to study the influence of the unobservables on the expected gains. With the data generating process of this simulation study it is clear both unobservables have an effect. It is thus essential to account for these two different sources of unobserved heterogeneity to get unbiased estimators of treatment effects parameters such as ATE or TT. Obviously, the choice of a domain of integration

<sup>25</sup>Recall that the estimator of Gautier and Le Penne (2011) is adaptive and thus achieves the minimax rate of convergence without having to choose a smoothing parameter.

is important when integration with respect to  $(\gamma, \theta)$  has to be carried out. We recommend defining the domain of integration as the set  $\{(\gamma, \theta) : \widehat{f}_{\Gamma, \Theta}(\gamma, \theta) > \tau\}$  or equivalently to calculate an integral against  $\widehat{f}_{\Gamma, \Theta}(\gamma, \theta) \mathbf{1}\{\widehat{f}_{\Gamma, \Theta}(\gamma, \theta) > \tau\}$ . In this simulation study, the box  $\mathcal{B}_1 = [-1.5, 3.5] \times [-3, 2]$  corresponds to the choice  $\tau \approx 0.021 \left\| \widehat{f}_{\Gamma, \Theta} \right\|_{\infty}$  while the box  $\mathcal{B}_2 = [-1.75, 3.75] \times [-3.25, 2.25]$  (which seems to perform better to estimate the ATE and TT) corresponds to the choice  $\tau \approx 0.001 \left\| \widehat{f}_{\Gamma, \Theta} \right\|_{\infty}$ . The domain  $\mathcal{B}_1$  contains 90% of the total mass while the domain  $\mathcal{B}_2$  contains 94% of the total mass<sup>26</sup>.

## 6. ALTERNATIVE APPROACH AND EXTENSIONS

**6.1. An Alternative Scaling of the Random Coefficients Binary Choice Model.** In this section, we present a different estimation approach based on the scaling in Ichimura and Thomson (1998), Gautier and Kitamura (2009) and Gautier and Le Penneç (2011). We do not condition on control variables for simplicity of the notations but it could be done exactly as it was done earlier. Equation (1.2) is of the form  $D = \mathbf{1}\{(V, Z', 1)(1, -\Gamma', -\Theta)' > 0\}$ . Because the scale of  $(1, -\Gamma', -\Theta)$  is not identified, instead of normalizing the first coordinate (the coefficient of  $V$ ) to be one, we can work with the vector  $\bar{\Gamma} = (1, -\Gamma', -\Theta) / \|(1, -\Gamma', -\Theta)\|$ <sup>27</sup> which is of norm 1. This yields more flexibility because a sufficient<sup>28</sup> condition for identification is that the support of  $\bar{\Gamma}$  belongs to an hemisphere (see Gautier and Kitamura (2009)) and it is not required that in the original scale one specific coefficient is positive. This is satisfied for example when a coefficient has a sign, but in that case the identity of the regressor which has a sign, and the sign itself, do not need to be known in advance.

**Remark 6.1.** The condition that the support of  $\bar{\Gamma}$  belongs to an hemisphere implies that, applying a rotation to the vector of instruments, one transformed instrument has a positive coefficient. This rotation does not have to be known. Recall that this condition is only sufficient for identification.

We also rescale the instruments so that  $S = (V, Z', 1)' / \|(V, Z', 1)\|$ . Both  $S$  and  $\bar{\Gamma}$  belong to the sphere  $\mathbb{S}^L$  of the Euclidian space  $\mathbb{R}^{L+1}$ .

We will now use the notation  $\sigma$  for the spherical measure on  $\mathbb{S}^L$ . The spaces  $L^p(\mathbb{S}^L)$  are the classical  $L^p$  spaces with respect to the measure  $\sigma$ . We denote by  $H^+ = \{s \in \mathbb{S}^L : s_{L+1} > 0\}$ . Odd, respectively even, functions are the closure in  $L^1(\mathbb{S}^L)$  of continuous functions such that  $\forall s \in$

<sup>26</sup>Indeed we have  $\int_{-4}^6 \int_{-5.5}^{4.5} \mathbb{E}_S \left[ \widehat{f}_{\Gamma, \Theta} \right] (\gamma, \theta) d\gamma d\theta \approx 1.014$  while  $\int_{\mathcal{B}_1} \mathbb{E}_S \left[ \widehat{f}_{\Gamma, \Theta} \right] (\gamma, \theta) d\gamma d\theta \approx 0.916$  and  $\int_{\mathcal{B}_2} \mathbb{E}_S \left[ \widehat{f}_{\Gamma, \Theta} \right] (\gamma, \theta) d\gamma d\theta \approx 0.952$ .

<sup>27</sup>We already used the notation  $\bar{\Gamma}$  in Section 3.4.1, we would like to warn the reader that here it is denoting a different quantity.

<sup>28</sup>Not a necessary condition.

$\mathbb{S}^L$ ,  $f(-s) = -f(s)$ , respectively  $\forall s \in \mathbb{S}^L$ ,  $f(-s) = f(s)$ . Each function in  $L^2(\mathbb{S}^L)$  is the orthogonal sum of its odd and even part. We denote by  $f^-$ , respectively  $f^+$ , the odd and even parts of a function  $f$  in  $L^1(\mathbb{S}^L)$ .

Under the above scaling, (1.2) becomes

$$(6.1) \quad D = \mathbf{1}\{S'\bar{\Gamma} > 0\}.$$

We make the following assumption.

**Assumption 6.1.** (A-1) The rescaled vector of instruments  $S$  has a density with respect to  $\sigma$  and

its support is the whole hemisphere  $\overline{H^+} = \{s \in \mathbb{S}^L : s_{L+1} \geq 0\}$  ;

(A-2)  $\bar{\Gamma}$  has a density  $f_{\bar{\Gamma}}$  with respect to  $\sigma$  which is defined point-wise and has support included in some hemisphere  $H = \{s \in \mathbb{S}^L : s' \mathbf{n} \geq 0\}$ , where  $\mathbf{n}$  is a vector of norm 1 that does not need to be known ;

(A-3) For the function  $\phi$  considered,  $\left(\overline{\mathbb{E}[\phi(Y_j)|\bar{\Gamma} = \cdot]} f_{\bar{\Gamma}}\right)^-$  belongs to  $L^2(\mathbb{S}^L)$  ;

(A-4)  $S_{\perp}(Y_0, \bar{\Gamma}')$  and  $S_{\perp}(Y_1, \bar{\Gamma}')$ .

Assumption (A-1) corresponds to full support of the instruments. It is stronger than Assumption 2.2. Assumption (A-2) is satisfied under the specification (1.1) where, in the original scale, the coefficient of  $V$  has a sign which is known. As explained, it is more general. Note that, in the generalized Roy model example, the random coefficients are cost factors and assuming that one coefficient as a sign is very credible. Several rescaling yielding an instruments being called  $V$  in (1.1) can be used. One should in theory be very cautious as some coefficients could have very small values for certain individuals yielding very large in absolute value coefficients in the new scale. One numerical advantage of the normalization of this section is that it avoids the possible arbitrariness of the choice of a regressor  $V$  and unstable division by numbers potentially close to 0 for some individuals. Assumption (A-4) corresponds to Assumption 2.1 (A-2).

**Theorem 6.1.** Under Assumption 6.1, for an arbitrary measurable function  $\phi$  such that  $E[|\phi(Y_0)| + |\phi(Y_1)|] < \infty$ , we obtain, for  $j = 1$  when  $\zeta(D) = D$  and  $j = 0$  when  $\zeta(D) = D - 1$ , for almost every  $\gamma$  in  $\mathbb{S}^L$ ,

$$(6.2) \quad \overline{\mathbb{E}[\phi(Y_j)|\bar{\Gamma} = \cdot]}(\gamma) f_{\bar{\Gamma}}(\gamma) = 2 \left(\overline{\mathbb{E}[\phi(Y_j)|\bar{\Gamma} = \cdot]} f_{\bar{\Gamma}}\right)^-(\gamma) \mathbf{1}\{f_{\bar{\Gamma}}(\gamma) > 0\}$$

where

$$(6.3) \quad \left(\overline{\mathbb{E}[\phi(Y_j)|\bar{\Gamma} = \cdot]} f_{\bar{\Gamma}}\right)^-(\gamma) = \mathcal{H}^{-1}(R_j),$$

$$\begin{aligned} \forall s \in H^+, R_j(s) &= \frac{1}{2} \mathbb{E}[\phi(Y_j)] - \mathbb{E}[\zeta(D)\phi(Y)|S = s] \\ \forall s \in -H^+, R_j(s) &= -R_j(-s), \end{aligned}$$

and for any point  $\tilde{s}$  on  $\partial H^+ = \overline{H^+} \setminus H^+$ ,

$$(6.4) \quad \mathbb{E}[\phi(Y_j)] = \lim_{s \rightarrow \tilde{s}, s \in H^+} \mathbb{E}[\zeta(D)\phi(Y)|S = s] + \lim_{s \rightarrow -\tilde{s}, s \in H^+} \mathbb{E}[\zeta(D)\phi(Y)|S = s]$$

The operator  $\mathcal{H}$  is the hemispherical transform. Let us recall a few of its properties (see, *e.g.*, Gautier and Kitamura (2009)). The operator  $\mathcal{H}$  is not injective in  $L^2(\mathbb{S}^L)$  but it is when restricted to  $L^2_{\text{odd}}(\mathbb{S}^L)$ , the closure in  $L^2(\mathbb{S}^L)$  of continuous and bounded odd functions. Also, the smoothing properties of  $\mathcal{H}$  together with the Sobolev embeddings imply that functions in  $\mathcal{H}(L^2_{\text{odd}}(\mathbb{S}^L))$  are continuous.

Because the right hand-side of (6.4) does not depend on  $\tilde{s}$ , an efficient estimator should take into account all these relations for all  $\tilde{s}$  on the boundary of  $H^+$ . The result (6.4) holds for our original model (1.1)-(1.2) when the instruments have full support. It is not specific to a particular scaling or operator<sup>29</sup>. The main reasons behind the existence of these formulas are: (1) the linear index structure and (2) the smoothing properties of the operator in the inverse problem formulation<sup>30</sup>. As a consequence of (6.4), we obtain the following corollary.

**Corollary 6.1.** Under Assumption 6.1, for an arbitrary function  $\phi$  such that  $E[|\phi(Y_0)| + |\phi(Y_1)|] < \infty$ , we obtain, for  $j = 1$  when  $\zeta(D) = D$  and  $j = 0$  when  $\zeta(D) = D - 1$ , for any sequence  $(z_0(N), \dots, z_{L-1}(N))_{N \in \mathbb{N}}$ , such that  $\lim_{N \rightarrow \infty} \|(z_0(N), \dots, z_{L-1}(N), 1)\| = \infty$ ,

$$\begin{aligned} \mathbb{E}[\phi(Y_j)] &= \lim_{N \rightarrow \infty} \left\{ \mathbb{E} \left[ \zeta(D)\phi(Y) \left| \frac{(V, Z', 1)}{\|(V, Z', 1)\|} = \frac{(z_0(N), \dots, z_{L-1}(N), 1)}{\|(z_0(N), \dots, z_{L-1}(N), 1)\|} \right. \right] \right. \\ &\quad \left. + \mathbb{E} \left[ \zeta(D)\phi(Y) \left| \frac{(V, Z', 1)}{\|(V, Z', 1)\|} = \frac{(-z_0(N), \dots, -z_{L-1}(N), 1)}{\|(-z_0(N), \dots, -z_{L-1}(N), 1)\|} \right. \right] \right\}; \end{aligned}$$

we obtain as well, for any  $l = 0, \dots, L - 1$ ,

$$(6.5) \quad \begin{aligned} \mathbb{E}[\phi(Y_j)] &= \lim_{z_l \rightarrow \infty} \mathbb{E} \left[ \zeta(D)\phi(Y) \left| \frac{(V, Z', 1)}{\|(V, Z', 1)\|} = \frac{(z_0, z_1, \dots, z_{L-1}, 1)}{\|(z_0, z_1, \dots, z_{L-1}, 1)\|} \right. \right] \\ &\quad + \lim_{z_l \rightarrow -\infty} \mathbb{E} \left[ \zeta(D)\phi(Y) \left| \frac{(V, Z', 1)}{\|(V, Z', 1)\|} = \frac{(z_0, z_1, \dots, z_{L-1}, 1)}{\|(z_0, z_1, \dots, z_{L-1}, 1)\|} \right. \right]. \end{aligned}$$

<sup>29</sup>The operator does not appear in the result and, when the instruments in (1.2) have full support, (1.2) can be transformed into (6.1) without loss of generality.

<sup>30</sup>In the proof we use the scaling of this section and the smoothing properties of the Hemispherical transform.



As a consequence, under the extra integrability condition (A-3), there exists a formula that identifies quantities related to the marginals, thus as well ATE or QTE, at infinity. Note that this is true by letting one of the instruments in  $Z$  going to infinity, even if this does not produce situations with “unselected samples”. One obtains for example, if  $\left(\overline{\mathbb{E}[Y_0|\overline{\Gamma} = \cdot]}f_{\overline{\Gamma}}\right)^-$  and  $\left(\overline{\mathbb{E}[Y_1|\overline{\Gamma} = \cdot]}f_{\overline{\Gamma}}\right)^-$  belong to  $L^2(\mathbb{S}^L)$ , that

$$(6.6) \quad \text{ATE} = \lim_{s \rightarrow \bar{s}, s \in H^+} \mathbb{E}[(2D - 1)Y|S = s] + \lim_{s \rightarrow -\bar{s}, s \in H^+} \mathbb{E}[(2D - 1)Y|S = s].$$

Given an estimator  $\mathbb{E}[\widehat{\phi}(Y_j)]$  of  $\mathbb{E}[\phi(Y_j)]$ , which is either obtained for large values of the instruments building on (6.4) or using the approach of Section 4, we can get the following estimator of  $\overline{\mathbb{E}[\phi(Y_j)|\overline{\Gamma} = \cdot]}(\gamma)f_{\overline{\Gamma}}$

$$(6.7) \quad \left(\overline{\widehat{\mathbb{E}[\phi(Y_j)|\overline{\Gamma} = \cdot]}f_{\overline{\Gamma}}}\right)^-(\gamma) = \frac{1}{|\mathbb{S}^L|} \sum_{p=0}^{T_N-1} \frac{\chi(2p+1, 2T_N)h(2p+1, L)}{\lambda(2p+1, L)C_{2p+1}^{\nu(L)}(1)} \left( \frac{1}{N} \sum_{i=1}^N \frac{\left(\mathbb{E}[\widehat{\phi}(Y_j)] - 2\zeta(d_i)\phi(y_i)\right) C_{2p+1}^{\nu(L)}(s'\gamma)}{\max(\hat{f}_S(s_i), m_N)} \right),$$

where  $|\mathbb{S}^L| = \frac{2\pi^{(L+1)/2}}{\Gamma((L+1)/2)}$  is the surface measure of  $\mathbb{S}^L$ ,  $h(n, L) = \frac{(2n+L-1)(n+L-1)!}{n!(L-1)!(n+L-1)}$ ,  $\nu(L) = (L - 1)/2$ ,  $\lambda(2p+1, L) = \frac{(-1)^p |\mathbb{S}^{L-1}| 1 \cdot 3 \cdots (2p-1)}{(L)(L+2)\cdots(L+2p)}$ ,  $\chi(n, T) = \psi(n/T)$  where  $\psi : [0, \infty) \rightarrow [0, \infty)$  is infinitely differentiable, nonincreasing, such that  $\psi(x) = 1$  if  $x \in [0, 1]$ ,  $0 \leq \psi(x) \leq 1$  if  $x \in [1, 2]$ ,  $\psi(x) = 0$  if  $x \geq 2$ , and  $C_n^{\nu}(\cdot)$  are the Gegenbauer polynomials. The Gegenbauer polynomials are given by

$$C_n^{\nu}(t) = \sum_{l=0}^{\lfloor n/2 \rfloor} \frac{(-1)^l (\nu)_{n-l}}{l!(n-2l)!} (2t)^{n-2l}, \quad \nu > -1/2, n \in \mathbb{N}$$

where  $(a)_0 = 1$  and for  $n$  in  $\mathbb{N} \setminus \{0\}$ ,  $(a)_n = a(a+1)\cdots(a+n-1) = \Gamma(a+n)/\Gamma(a)$ .  $T_N$  is the smoothing parameter,  $m_N$  a trimming factor and  $\hat{f}_S$  an estimator of the density of  $S$ . This estimator is in the same spirit as in the same spirit as the estimator in Gautier and Kitamura (2009) (see the reference for more details). Without plug-in, the method of Gautier and Le Pennec (2011) which is a powerful completely data driven adaptive method could also be used.

One main drawback of the approach of this section is that it relies on plug-in estimators of  $\mathbb{E}[\widehat{\phi}(Y_j)]$ . Using plug-in estimators from Section 4 is not very satisfactory because the method of this section is no longer completely alternative. Using (6.4) to obtain the plug-in is very inefficient as it only uses the large values of the instruments<sup>31</sup> and relies critically on the integrability condition (A-3).

The rest of the paper does not rely on such formulas involving only values at infinity of the instruments. On the contrary, the estimators that we consider use all the observations and involve trimming of values of the instruments in the tails. In order to obtain a much larger class of treatment

<sup>31</sup>Which in practice are often mis-measured or outliers.

effect parameters one calculates weighted integrals of UCATE or UCDITE on a domain  $\mathcal{B}_N$  which again acts as trimming of values of the instruments in the tails. Apart from being mostly inapplicable for estimation, formulas of the type of (6.4) do not yield the roots UCATE or UCDITE. This is important because we have seen that we need to control for unobserved heterogeneity to justify an assumption such as Assumption 3.2 and obtain treatment effects that depend on the whole distribution of potential outcomes.

**6.2. The Case of Binary Instruments.** We have seen that Assumption 2.1 (A-1) allows to some extent, when  $L \geq 2$ , cases where  $V$  is discrete. One needs the strong support condition (A-2.2). In this section, we consider the case where instruments other than  $V$  are binary. We replace (1.2) by

$$(6.8) \quad D = \mathbf{1} \{V - \alpha B - \Gamma' Z - \Theta > 0\}$$

where  $B$  is a binary instrument and  $(\Gamma', \Theta, \alpha)$  is a vector of random coefficients of dimension  $L + 1$ . We rewrite equation (6.8) in the form

$$(6.9) \quad D = \mathbf{1} \left\{ \tilde{S}'((\Gamma', \Theta + \alpha B)') < \tilde{V} \right\}$$

where  $\tilde{S}$  and  $\tilde{V}$  are defined in Section 2. We make the following assumption.

**Assumption 6.2.** (A-1) The conditional distribution of  $(\tilde{S}', \tilde{V}, \Gamma', \Theta, \alpha)$  given  $X = x$  is absolutely continuous with respect to the product of the spherical measure on  $\mathbb{S}^{L-1}$  and the Lebesgue measure on  $\mathbb{R}^{L+2}$  for almost every  $x$  in  $\text{supp}(X)$  ;

(A-2)  $(V, Z, B) \perp (Y_0, \Gamma', \Theta, \alpha) | X$  and  $(V, Z, B) \perp (Y_1, \Gamma', \Theta, \alpha) | X$  ;

(A-3)  $0 < \mathbb{P}(D = 1 | X) < 1$  a.s. ;

(A-4)  $X_0 = X_1$  a.s. ;

(A-5) (case 1) for every  $x \in \text{supp}(X)$ ,  $\text{supp}(f_{\tilde{S}|X, B=1}(\cdot | x)) = \overline{H^+}$  and for every  $s \in \text{Int}(H^+)$ ,

$$\text{supp}(f_{\tilde{V}|\tilde{S}, X, B=1}(\cdot | s, x)) \supset \left[ \inf_{\gamma \in \text{supp}(f_{\Gamma, \Theta + \alpha | X}(\cdot | x))} s' \gamma, \sup_{\gamma \in \text{supp}(f_{\Gamma, \Theta + \alpha | X}(\cdot | x))} s' \gamma \right].$$

or (case 2) for every  $x \in \text{supp}(X)$ ,  $\text{supp}(f_{\tilde{S}|X, B=0}(\cdot | x)) = \overline{H^+}$  and for every  $s \in \text{Int}(H^+)$ ,

$$\text{supp}(f_{\tilde{V}|\tilde{S}, X, B=0}(\cdot | s, x)) \supset \left[ \inf_{\gamma \in \text{supp}(f_{\Gamma, \Theta | X}(\cdot | x))} s' \gamma, \sup_{\gamma \in \text{supp}(f_{\Gamma, \Theta | X}(\cdot | x))} s' \gamma \right].$$

**Theorem 6.2.** Consider an arbitrary function  $\phi$  such that  $\mathbb{E}[|\phi(Y_0)| + |\phi(Y_1)|] < \infty$ . Let  $L \geq 2$ , and make Assumption 6.2. Then, in case 1, the following statements hold, almost surely in  $x$  in  $\text{supp}(X)$ ,

$$(6.10) \quad f_{\Gamma, \Theta + \alpha | X}(\cdot | x) = R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ D \mid \left( \tilde{S}, \tilde{V} \right) = \cdot, B = 1, X = x \right]} \right]$$

$$(6.11) \quad \overline{\mathbb{E} [\phi(Y_1) | (\Gamma, \Theta + \alpha) = \cdot, X = x]} f_{\Gamma, \Theta + \alpha | X}(\cdot | x) = R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ \phi(Y) D \mid \left( \tilde{S}, \tilde{V} \right) = \cdot, B = 1, X = x \right]} \right]$$

$$(6.12) \quad \overline{\mathbb{E} [\phi(Y_0) | (\Gamma, \Theta + \alpha) = \cdot, X = x]} f_{\Gamma, \Theta + \alpha | X}(\cdot | x) = R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ \phi(Y) (D - 1) \mid \left( \tilde{S}, \tilde{V} \right) = \cdot, B = 1, X = x \right]} \right]$$

while in case 2, the following statements hold, almost surely in  $x$  in  $\text{supp}(X)$ ,

$$(6.13) \quad f_{\Gamma, \Theta | X}(\cdot | x) = R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ D \mid \left( \tilde{S}, \tilde{V} \right) = \cdot, B = 0, X = x \right]} \right]$$

$$(6.14) \quad \overline{\mathbb{E} [\phi(Y_1) | (\Gamma, \Theta) = \cdot, X = x]} f_{\Gamma, \Theta | X}(\cdot | x) = R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ \phi(Y) D \mid \left( \tilde{S}, \tilde{V} \right) = \cdot, B = 0, X = x \right]} \right]$$

$$(6.15) \quad \overline{\mathbb{E} [\phi(Y_0) | (\Gamma, \Theta) = \cdot, X = x]} f_{\Gamma, \Theta | X}(\cdot | x) = R^{-1} \left[ \overline{\partial_v \mathbb{E} \left[ \phi(Y) (D - 1) \mid \left( \tilde{S}, \tilde{V} \right) = \cdot, B = 0, X = x \right]} \right].$$

This result is straightforward to obtain using the same arguments as those that allowed to prove Theorem 6.2. Estimators could be obtained like in Section 4.1. Theorem 3.1 yields that, in case 1,  $\text{UCATE}(\gamma, \alpha + \theta, x)$  and  $f_{\Gamma, \Theta + \alpha}$  are identified and estimable, thus all treatment effect parameters that depend on averages can be estimated. In case 2,  $\text{UCATE}(\gamma, \theta, x)$  and  $f_{\Gamma, \Theta}$  are also estimable and allow to estimate all treatment effect parameters that depend on averages. To obtain treatment effect parameters that depend on the distribution of potential outcomes one needs, in case 1, to replace Assumption 3.2 by

**Assumption 6.3.**  $Y_0 \perp \Delta \mid \Gamma, \Theta + \alpha, X$ .

While in case 2 one can simply rely on Assumption 3.2. Note that when Assumption 6.2 holds both in case 1 and 2. This yields 2 formulas to identify the various treatment effects that depend on averages. When as well both Assumption 3.2 and Assumption hold then we also obtain 2 formulas to identify treatment effects that depend on the distribution of potential outcomes. It is the possible to combine the 2 estimators, built on the sub-samples where  $b_i = 1$  and  $b_i = 0$  respectively, to make a more efficient use of the data.

**Remark 6.2.** Note that we do not need to estimate the full joint distribution  $f_{\Gamma, \Theta, \alpha|X}$  nor  $\mathbb{E}[\Delta|\Gamma, \Theta, \alpha, X]$  or  $f_{\Delta|\Gamma, \Theta, \alpha, X}$  to obtain the various treatment effect parameters. Estimating these parameters requires more assumptions. For example assuming that  $\Theta \perp \alpha|\Gamma, X$  is enough to identify  $f_{\Gamma, \Theta, \alpha|X}$  from  $f_{\Gamma, \Theta + \alpha|X}$  and  $f_{\Gamma, \Theta|X}$  by conditional deconvolution. This last assumption is of the same nature as Assumption 3.2 and allows to identify the joint distribution of random coefficients in a linear model with a binary.

## 7. APPENDIX

7.0.1. *Proof of Proposition 4.1.* Because  $K_T$  and  $\tilde{K}_T$  belong to  $\mathcal{S}(\mathbb{R})$ , due to Assumption 4.1 (ii), for the function  $\phi$  considered, for almost every  $s$  in  $H^+$ ,  $v \mapsto \mathbb{E}[\phi(Y)\zeta(D)|(\tilde{S}, \tilde{V}) = (s, v)]\tilde{K}_T(v)$  and  $v \mapsto \partial_v \mathbb{E}[\phi(Y)\zeta(D)|(\tilde{S}, \tilde{V}) = (s, v)]K_T(v)$  are in  $L^1(\mathbb{R})$  and  $\lim_{|v| \rightarrow \infty} \mathbb{E}[\phi(Y)\zeta(D)|(\tilde{S}, \tilde{V}) = (s, v)]K_T(v) = 0$ . This yields

$$\begin{aligned}
& A_T \left[ \overline{\partial_v \mathbb{E} \left[ \phi(Y)\zeta(D) | (\tilde{S}, \tilde{V}) = (\cdot), X = x \right]} \right] (\gamma) \\
&= \int_{\text{supp}(f_{\tilde{S}, \tilde{V}|X}(\cdot|x))} \tilde{K}_T(s'\gamma - u) \mathbb{E} \left[ \phi(Y)\zeta(D) | (\tilde{S}, \tilde{V}) = (s, u), X = x \right] dud\sigma(s) \quad (\text{by integration by parts}) \\
&= \int_{\text{supp}(f_{\tilde{S}, \tilde{V}|X}(\cdot|x))} \tilde{K}_T(s'\gamma - u) \mathbb{E} \left[ \phi(Y)\zeta(D) | (\tilde{S}, \tilde{V}) = (s, u), X = x \right] \frac{f_{\tilde{S}, \tilde{V}|X}(s, u|x)}{f_{\tilde{S}, \tilde{V}|X}(s, u|x)} dud\sigma(s) \\
&= \int_{\text{supp}(f_{\tilde{S}, \tilde{V}|X}(\cdot|x))} \mathbb{E} \left[ \frac{\tilde{K}_T(s'\gamma - u)\phi(Y)\zeta(D)}{f_{\tilde{S}, \tilde{V}|X}(s, u|x)} \middle| (\tilde{S}, \tilde{V}) = (s, u), X = x \right] f_{\tilde{S}, \tilde{V}|X}(s, u|x) dud\sigma(s) \\
&= \mathbb{E} \left[ \frac{\tilde{K}_T(\tilde{S}^T \gamma - U)\phi(Y)\zeta(D)}{f_{\tilde{S}, \tilde{V}|X}(\tilde{S}, \tilde{V}|x)} \middle| X = x \right] \quad (\text{by the law of iterated conditional expectations})
\end{aligned}$$

Q.E.D.

7.0.2. *Proof of Proposition 4.2.* We use the notations

$$\begin{aligned}
g_{m, \tau}^I(\gamma) &= \frac{1}{N} \sum_{i=1}^N \frac{\tilde{K}_{T_N}(\tilde{s}'_i \gamma - \tilde{v}_i) T_{\tau_N}(\phi(y_i)) \zeta(d_i)}{\max(f_{\tilde{S}, \tilde{V}}(\tilde{s}_i, \tilde{v}_i), m_N)}, \\
g_{\tau}^I(\gamma) &= \frac{1}{N} \sum_{i=1}^N \frac{\tilde{K}_{T_N}(\tilde{s}'_i \gamma - \tilde{v}_i) T_{\tau_N}(\phi(y_i)) \zeta(d_i)}{f_{\tilde{S}, \tilde{V}}(\tilde{s}_i, \tilde{v}_i)}, \\
g^I(\gamma) &= \frac{1}{N} \sum_{i=1}^N \frac{\tilde{K}_{T_N}(\tilde{s}'_i \gamma - \tilde{v}_i) \phi(y_i) \zeta(d_i)}{f_{\tilde{S}, \tilde{V}}(\tilde{s}_i, \tilde{v}_i)},
\end{aligned}$$

where the superscript  $I$  stands for ideal (this is because we replace the estimator of the density in the denominator by the true density). For two sequences of positive numbers  $(a_n)_{n \in \mathbb{N}}$  and  $(b_n)_{n \in \mathbb{N}}$ ,

we write  $a_n \lesssim b_n$  when there exists  $M$  positive such that  $a_n \leq Mb_n$  and  $a_n \asymp b_n$  when  $a_n \lesssim b_n$  and  $b_n \lesssim a_n$ .

Let us start by stating a few results on  $\tilde{K}_T$ . Recall that

$$\tilde{K}_T(u) = \frac{2}{(2\pi)^L} \int_0^\infty \sin(-tu)t|t|^{L-1}\psi\left(\frac{t}{T}\right) dt$$

therefore by the change of variables

$$\tilde{K}_T(u) = \frac{2T^{L+1}}{(2\pi)^L} \int_0^\infty \sin(-Ttu)t|t|^{L-1}\psi(t)dt$$

thus

$$\left| \tilde{K}_T(u) \right| \leq \frac{2T^{L+1}}{(2\pi)^L} \int_0^\infty t^L \psi(t) dt$$

and  $\int_0^\infty t^L \psi(t) dt$  is a constant independent of  $T$  because  $\psi \in \mathcal{S}(\mathbb{R}^L)$ , therefore

$$(7.1) \quad \left| \tilde{K}_T \right|_\infty \lesssim T^{L+1}.$$

Similarly we can show that

$$(7.2) \quad \left| \tilde{K}'_T \right|_\infty \lesssim T^{L+2}$$

which implies that  $\forall (u, v) \in \mathbb{R}^2$ ,

$$\left| \tilde{K}_T(u) - \tilde{K}_T(v) \right|_\infty \lesssim T^{L+2}|u - v|$$

which in turns yields

$$\left| |\tilde{K}_T(u)| - |\tilde{K}_T(v)| \right|_\infty \lesssim T^{L+2}|u - v|$$

and  $\forall (s, v) \in \mathbb{S}^{L-1} \times \mathbb{R}$ ,

$$(7.3) \quad \left| \tilde{K}_T(s'\gamma - v) - \tilde{K}_T(s'\bar{\gamma} - v) \right| \lesssim T^{L+2}|\gamma - \bar{\gamma}|$$

$$(7.4) \quad \left| |\tilde{K}_T(s'\gamma - v)| - |\tilde{K}_T(s'\bar{\gamma} - v)| \right| \lesssim T^{L+2}|\gamma - \bar{\gamma}|$$

where on the right hand-side of (7.3) and (7.4)  $|\cdot|$  is the Euclidian norm in  $\mathbb{R}^L$ .

Also, because

$$\tilde{K}_T(u) = \frac{1}{(2\pi)^L} \int_{-\infty}^\infty e^{-iut}t|t|^{L-1}\psi\left(\frac{t}{T}\right) dt,$$

$$\begin{aligned} \left| \tilde{K}_T \right|_2 &= \left\{ \int_{-\infty}^\infty t^{2L} \psi^2\left(\frac{t}{T}\right) dt \right\}^{1/2} \\ &= T^{(2L+1)/2} \left\{ \int_{-\infty}^\infty t^{2L} \psi^2(t) dt \right\}^{1/2} \end{aligned}$$

$$(7.5) \quad \lesssim T^{(2L+1)/2}.$$

We now rely on the decomposition

$$\begin{aligned} \hat{g} - g &= (\hat{g} - g_{m,\tau}^I) + (g_{m,\tau}^I - \mathbb{E}[g_{m,\tau}^I]) + (\mathbb{E}[g_{m,\tau}^I] - \mathbb{E}[g_\tau^I]) + (\mathbb{E}[g_\tau^I] - \mathbb{E}[g^I]) + (\mathbb{E}[g^I] - g) \\ &:= S_p + S_e + B_t + B_{\text{trunc}} + B_a \end{aligned}$$

where the expectation is with respect to  $(\tilde{S}_i, \tilde{V}_i)_{i=1}^N$ <sup>32</sup>. The contribution  $S_p$  corresponds to the stochastic component due to plug-in,  $S_e$  to the stochastic component of the infeasible estimator  $g_{m,\tau}^I$ ,  $B_t$  to the trimming bias,  $B_{\text{trunc}}$  to the bias due to truncation and  $B_a$  to the approximation bias.

Let us study first the contribution of the term  $S_p$ . The following upper bounds hold

$$\begin{aligned} \|S_p\|_\infty &= \left\| \frac{1}{N} \sum_{i=1}^N \frac{\tilde{K}_{T_N}(\tilde{s}'_i \cdot -\tilde{v}_i) T_{\tau_N}(\phi(y_i)) \varsigma(d_i)}{\max(f_{\tilde{S}, \tilde{V}}(\tilde{s}_i, \tilde{v}_i), m_N)} \left( \frac{\max(f_{\tilde{S}, \tilde{V}}(\tilde{s}_i, \tilde{v}_i), m_N)}{\max(\widehat{f_{\tilde{S}, \tilde{V}}}(\tilde{s}_i, \tilde{v}_i), m_N)} - 1 \right) \mathbf{1}\{\cdot \in \mathcal{B}_N\} \right\|_\infty \\ &\leq \min(\tau_N, \|\phi\|_\infty) \left\| \frac{1}{N} \sum_{i=1}^N \frac{|\tilde{K}_{T_N}(\tilde{s}'_i \cdot -\tilde{v}_i)| \mathbf{1}\{\cdot \in \mathcal{B}_N\}}{\max(f_{\tilde{S}, \tilde{V}}(\tilde{s}_i, \tilde{v}_i), m_N)} \right\|_\infty \max_{i=1, \dots, N} \left| \frac{\max(f_{\tilde{S}, \tilde{V}}(\tilde{s}_i, \tilde{v}_i), m_N)}{\max(\widehat{f_{\tilde{S}, \tilde{V}}}(\tilde{s}_i, \tilde{v}_i), m_N)} - 1 \right| \\ &\leq m_N^{-1} \min(\tau_N, \|\phi\|_\infty) \left\| \frac{1}{N} \sum_{i=1}^N \frac{|\tilde{K}_{T_N}(\tilde{s}'_i \cdot -\tilde{v}_i)| \mathbf{1}\{\cdot \in \mathcal{B}_N\}}{\max(f_{\tilde{S}, \tilde{V}}(\tilde{s}_i, \tilde{v}_i), m_N)} \right\|_\infty \max_{i=1, \dots, N} |f_{\tilde{S}, \tilde{V}}(\tilde{s}_i, \tilde{v}_i) - \widehat{f_{\tilde{S}, \tilde{V}}}(\tilde{s}_i, \tilde{v}_i)| \\ &\leq m_N^{-1} \min(\tau_N, \|\phi\|_\infty) (\|T_1\|_\infty + \|T_2 \mathbf{1}\{\mathcal{B}_N\}\|_\infty) \max_{i=1, \dots, N} |f_{\tilde{S}, \tilde{V}}(\tilde{s}_i, \tilde{v}_i) - \widehat{f_{\tilde{S}, \tilde{V}}}(\tilde{s}_i, \tilde{v}_i)| \end{aligned}$$

where

$$\begin{aligned} T_1(\gamma) &= \mathbb{E} \left[ \frac{|\tilde{K}_{T_N}(\tilde{S}'\gamma - \tilde{V})|}{\max(f_{\tilde{S}, \tilde{V}}(\tilde{S}, \tilde{V}), m_N)} \right] \\ T_2(\gamma) &= \frac{1}{N} \sum_{i=1}^N \left\{ \frac{|\tilde{K}_{T_N}(\tilde{s}'_i \gamma - \tilde{v}_i)|}{\max(f_{\tilde{S}, \tilde{V}}(\tilde{s}_i, \tilde{v}_i), m_N)} - \mathbb{E} \left[ \frac{|\tilde{K}_{T_N}(\tilde{S}'\gamma - \tilde{V})|}{\max(f_{\tilde{S}, \tilde{V}}(\tilde{S}, \tilde{V}), m_N)} \right] \right\} \end{aligned}$$

We just have to consider the term  $\|T_2 \mathbf{1}\{\mathcal{B}_N\}\|_\infty$ . We cover  $\mathcal{B}_N$  by  $\mathfrak{N}(N, L)$  Euclidian balls  $(B_i)_{i=1}^{\mathfrak{N}(N, L)}$  of centers  $(\tilde{\gamma}_i)_{i=1}^{\mathfrak{N}(N, L)}$  and radius  $R(N, L)$ . Because  $\mathcal{B}_N$  is compact we have  $\mathfrak{N}(N, L) \asymp d(\mathcal{B}_N)^L R(N, L)^{-L}$ .

For  $M(\alpha)$  positive and an appropriately chosen sequence  $(v_N)$  to be defined later

$$(7.6) \quad \mathbb{P}(v_N \|T_2 \mathbf{1}\{\mathcal{B}_N\}\|_\infty \geq M(\alpha))$$

<sup>32</sup>Thus we do not integrate against the distribution of  $\tilde{\Gamma}$ .

$$\begin{aligned} &\leq \mathbb{P} \left( \bigcup_{i=1, \dots, \mathfrak{N}(N, L)} \{v_N |T_2(\bar{\gamma}_i)| \geq M(\alpha)/2\} \right) \\ &\quad + \mathbb{P} \left( \exists i \in \{1, \dots, \mathfrak{N}(N, L)\} : v_N \sup_{\gamma \in B_i} |T_2(\gamma) - T_2(\bar{\gamma}_i)| \geq M(\alpha)/2 \right). \end{aligned}$$

(7.7)

By taking  $R(N, L) \asymp m_N v_N^{-1} T_N^{-(L+2)} M(\alpha)$  for a well chosen constant, the first term on the right hand-side is equal to zero. This follows from the fact that  $T_2$  is Lipschitz with a constant proportional to  $m_N^{-1} T_N^{-(L+2)}$ . This is a consequence of (7.4). For such a choice of  $R(N, L)$ ,

$$(7.8) \quad \mathbb{P}(v_N \|T_2\|_\infty \geq M(\alpha)) \leq \mathfrak{N}(N, L) \sup_{i=1, \dots, \mathfrak{N}(N, L)} \mathbb{P}(v_N |T_2(\bar{\gamma}_i)| \geq M(\alpha)/2).$$

Now

$$\begin{aligned} &\mathbb{P}(v_N |T_2(\bar{\gamma}_i)| \geq M(\alpha)/2) \\ &= \mathbb{P} \left( \left| \sum_{j=1}^N \frac{|\tilde{K}_{T_N}(\tilde{s}'_j \bar{\gamma}_i - \tilde{v}_j)|}{\max(f_{\tilde{S}, \tilde{V}}(\tilde{s}_i, \tilde{v}_i), m_N) m_N^{-1} T_N^{L+1}} - \mathbb{E} \left[ \frac{|\tilde{K}_{T_N}(\tilde{S}' \bar{\gamma}_i - \tilde{V})|}{\max(f_{\tilde{S}, \tilde{V}}(\tilde{S}, \tilde{V}), m_N) m_N^{-1} T_N^{L+1}} \right] \right| \geq t \right) \\ (7.9) \quad &\leq 2 \exp \left\{ -\frac{1}{2} \left( \frac{t^2}{\omega + Lt/3} \right) \right\} \quad (\text{Bernstein inequality}) \end{aligned}$$

where

$$\begin{aligned} t &= T_N^{-(L+1)} v_N^{-1} m_N N M(\alpha)/2 \\ \omega &\geq \sum_{j=1}^N \text{var} \left( \frac{|\tilde{K}_{T_N}(\tilde{s}'_j \bar{\gamma}_i - \tilde{v}_j)|}{m_N^{-1} T_N^{L+1}} \right) \\ L &\geq \sup_{(s, v) \in \text{supp}(\tilde{S}, \tilde{V})} \left| \frac{\tilde{K}_{T_N}(s' \bar{\gamma}_i - v)}{\max(f_{\tilde{S}, \tilde{V}}(\tilde{S}, \tilde{V}), m_N) m_N^{-1} T_N^{L+1}} \right| \asymp 1 \quad (\text{using (7.1)}). \end{aligned}$$

As

$$\begin{aligned} \sum_{j=1}^N \text{var} \left( \frac{|\tilde{K}_{T_N}(\tilde{S}'_j \bar{\gamma}_i - \tilde{V}_j)|}{\max(f_{\tilde{S}, \tilde{V}}(\tilde{S}, \tilde{V}), m_N) m_N^{-1} T_N^{L+1}} \right) &\leq \frac{m_N^2}{T_N^{2(L+1)}} \sum_{j=1}^N \mathbb{E} \left[ \left( \frac{|\tilde{K}_{T_N}(\tilde{S}' \bar{\gamma}_i - \tilde{V})|}{\max(f_{\tilde{S}, \tilde{V}}(\tilde{S}, \tilde{V}), m_N)} \right)^2 \right] \\ &\leq \frac{m_N N T_N^{2L+1}}{T_N^{2(L+1)}} \quad (\text{Due to (7.5)}) \end{aligned}$$

we shall take  $\omega = m_N N T_N^{2L+1} T_N^{-2(L+1)}$ .

Now choose  $v_N$  such that  $t \asymp M(\alpha) \sqrt{\omega \log(N)}$ . Thus  $\omega$  is the leading term in the denominator of the exponent in (7.9). The corresponding  $v_N$  is

$$\begin{aligned} v_N &\asymp (\log N)^{-1/2} \omega^{-1/2} T_N^{-(L+1)} m_N N \\ &\asymp \left( \frac{N}{\log N} \right)^{1/2} m_N^{1/2} T_N^{-(L+1/2)}. \end{aligned}$$

For these choices of the parameters

$$(7.10) \quad \frac{t^2}{\omega + Lt/3} \asymp (\log N) M(\alpha)^2$$

and

$$R(N, L) \asymp \left( \frac{N}{\log N} \right)^{-1/2} m_N^{1/2} T_N^{-3/2} M(\alpha).$$

Due to (7.5) and because by assumption  $\log(T_N^3/m_N) + L \log(d(\mathcal{B}_N)) \leq \alpha$ , we obtain

$$(7.11) \quad \mathfrak{R}(N, L) \asymp d(\mathcal{B}_N)^L R(N, L)^{-L} = \exp((\alpha + L/2) \log N + o(\log N)).$$

Equations (7.8), (7.9), (7.10) and (7.11) imply that, for a positive constants  $C$  and  $C_2$ ,

$$(7.12) \quad \mathbb{P} \left( \left( \frac{N}{\log N} \right)^{1/2} T_N^{-(L+1/2)} m_N^{1/2} \|T_2 \mathbf{1}\{\mathcal{B}_N\}\|_\infty \geq M(\alpha) \right) \leq C \exp \{ (\log N) ((\alpha + L/2) - C_2 M(\alpha)^2) \}$$

holds. For a large enough  $M(\alpha)$ ,  $(\alpha + L/2) - C_2 M(\alpha)^2 < -1$  which implies summability of the left hand-side in (7.12), hence by the first Borel-Cantelli lemma for  $M(\alpha)$  large enough with probability one

$$\overline{\lim}_{N \rightarrow \infty} \left( \frac{N}{\log N} \right)^{1/2} T_N^{-(L+1/2)} m_N^{1/2} \|T_2\|_\infty < M(\alpha).$$

In summary, we have obtained that for some constant  $M_{IV}$  and  $M(\alpha)$ , with probability one, for every  $\epsilon$  positive, there exists  $N$  large enough such that

$$\begin{aligned} \|S_p \mathbf{1}\{\mathcal{B}_N\}\|_\infty &\leq (M_{IV} + \epsilon) \min(\tau_N, \|\phi\|_\infty) r_{IV, N} m_N^{-1} \left\{ \left\| \mathbb{E} \left[ \frac{|\tilde{K}_{T_N}(\tilde{S}'\gamma - \tilde{V})|}{\max(f_{\tilde{S}, \tilde{V}}(\tilde{S}, \tilde{V}), m_N)} \right] \right\|_\infty \right. \\ &\quad \left. + m_N^{-1/2} (M(\alpha) + \epsilon) \left( \frac{N}{\log N} \right)^{-1/2} T_N^{L+1/2} \right\}. \end{aligned}$$

For the same reason, on the same event of probability 1, for every  $\epsilon$  positive, there exists  $N$  large enough such that

$$\|S_e \mathbf{1}\{\mathcal{B}_N\}\|_\infty \leq m_N^{-1/2} (M(\alpha) + \epsilon) \min(\tau_N, \|\phi\|_\infty) \left( \frac{N}{\log N} \right)^{-1/2} T_N^{L+1/2}.$$



Consider now the bias term induced by trimming, evaluated at a point  $\gamma$ ,

$$\begin{aligned} B_t(\gamma) &= \mathbb{E} \left[ \frac{\tilde{K}_{T_N}(\tilde{S}'\gamma - \tilde{V}) \min(\phi(Y), \tau_N) \varsigma(D)}{f_{\tilde{S}, \tilde{V}}(\tilde{S}, \tilde{V})} \left( \frac{f_{\tilde{S}, \tilde{V}}(\tilde{S}, \tilde{V})}{\max(f_{\tilde{S}, \tilde{V}}(\tilde{S}, \tilde{V}), m_N)} - 1 \right) \right] \\ &= \int_{\{(s,v): f_{\tilde{S}, \tilde{V}}(s,v) < m_N\}} \mathbb{E} \left[ \min(\phi(Y), \tau_N) \varsigma(D) | \tilde{S} = s, \tilde{V} = v \right] \tilde{K}_{T_N}(s'\gamma - v) \left( f_{\tilde{S}, \tilde{V}}(s,v) m_N^{-1} - 1 \right) d\sigma(s) dv. \end{aligned}$$

This yields the following upper bound

$$|B_t(\gamma)| \leq \min(\tau_N, \|\phi\|_\infty) \int_{\{(s,v): f_{\tilde{S}, \tilde{V}}(s,v) < m_N\}} \left| \tilde{K}_{T_N}(s'\gamma - v) \right| d\sigma(s) dv.$$

Consider now the truncation bias  $B_{\text{trunc}}$ . We obtain

$$|B_{\text{trunc}}| \leq \mathbb{E} \left[ \left| \tilde{K}_{T_N}(\tilde{S}^T \gamma - \tilde{V}) T_{\tau_N}(\phi(Y)) \zeta(D) \right| \right]$$

which allows to conclude using (7.1) with an explicit constant.

The upper bound for the approximation bias  $B_a$  is obtained as follows. Note that, for  $x$  in  $\mathbb{R}^L$ ,

$$\begin{aligned} \mathbb{E}[g^I](x) &= \mathcal{F}^{-1} \left[ \psi \left( \frac{\cdot}{T} \right) \mathcal{F}[g](\cdot) \right] (x) \\ &= \rho_T * g(x) \end{aligned}$$

where  $*$  is the usual convolution and

$$\begin{aligned} \rho_T(x) &= \frac{1}{(2\pi)^L} \int_{-\infty}^{\infty} e^{-i\xi'x} \psi \left( \frac{\xi}{T} \right) d\xi \\ (7.13) \quad &= T^L \rho_1 \left( \frac{x}{T} \right). \end{aligned}$$

The collection  $(\rho_T)_{T>0}$  is an approximate identity because of (7.13) and  $\int_{-\infty}^{\infty} \rho_T(x) dx = 1 (= \psi(0))$ .

The rest of the argument is classical and is based on

$$g - \mathbb{E}[g^I](x) = \int_{\mathbb{R}^L} (g(x) - g(x-y)) \rho_T(y) dy.$$

Let us do the argument for  $s = 1$  and  $s = 2$  only for simplicity of the notations.

**Case where  $s = 1$ .**

The inequalities

$$\begin{aligned} \|g - \mathbb{E}[g^I]\|_\infty &\leq L_g \int_{\mathbb{R}^L} |y| |\rho_T(y)| dy \\ &\leq \|g\|_{1,\infty} T^{-1} \int_{\mathbb{R}^L} |y| |\rho_1(y)| dy \end{aligned}$$

hold with  $L_g$  the Lipschitz constant of  $g$  which is itself upper bounded by  $\|g\|_{1,\infty}$ . The last integral is finite because  $\rho_1$  is in  $\mathcal{S}(\mathbb{R}^L)$ . Indeed  $\rho_1$  is the Fourier transform of a function in  $\mathcal{S}(\mathbb{R}^L)$ .

**Case where  $s = 2$ .**

Denoting by  $Dg(x).y$  the differential of  $g$  at  $x$  applied to  $y$ , because  $\psi$  and thus  $\rho_1$  is symmetric,

$$\|g - \mathbb{E}[g^I]\|_\infty = \int_{\mathbb{R}^L} (g(x) - g(x - y) - Dg(x).y)\rho_T(y)dy.$$

This yields

$$\begin{aligned} \|g - \mathbb{E}[g^I]\|_\infty &= \int_{\mathbb{R}^L} \int_0^1 (Dg(x - \lambda y) - Dg(x)).y d\lambda \rho_T(y) dy \\ &\leq \|g\|_{2,\infty} \int_{\mathbb{R}^L} |y|^2 |\rho_T(y)| dy \\ &\leq \|g\|_{2,\infty} T^{-2} \int_{\mathbb{R}^L} |y|^2 |\rho_1(y)| dy \end{aligned}$$

where again  $\int_{\mathbb{R}^L} |y|^2 |\rho_1(y)| dy < \infty$  because  $\rho_1$  is in  $\mathcal{S}(\mathbb{R}^L)$ .

The upper bounds on the bias due to truncation follow from the expression of the difference between the two expectations when Assumption 2.1 holds.

Q.E.D.

7.0.3. *Proof of Proposition 4.3.* The proof of this result is almost the same as the proof of Proposition 4.2. We will thus only stress the differences. We will use the notation

$$\|f - g\|_\infty := \sup_{t \in [-R_N^{\max}, R_N^{\max}], \gamma \in \mathcal{B}_N} |(f - g)(t, \gamma)|.$$

We start by observing that

$$\|f - g\|_\infty \leq \|\Re(f) - \Re(g)\|_\infty + \|\Im(f) - \Im(g)\|_\infty$$

this yields the factor 2 in the upper bound of the proposition. Then it is easy to check that we obtain the same upper bounds for both the error on the estimation of the real part and the error on the estimation of the imaginary part. For both, all the terms can be bounded like in the proof of Proposition 4.2 (taking  $\tau_N = 1$  and noting that  $\|\phi\|_\infty = 1$  as here  $\phi$  is either cos or sin) besides the term  $S_a$ , the stochastic component of the infeasible estimator  $g_{m,\tau}^I$ .

We shall cover  $\mathcal{B}_N$  by  $\overline{\mathfrak{N}}(N, L)$  balls  $(B_i)_{i=1}^{\overline{\mathfrak{N}}(N, L)}$  of centers  $(\bar{\gamma}_i, \bar{t}_i)_{i=1}^{\overline{\mathfrak{N}}(N, L)}$  and radius  $R(N, L)$  (will be the same as in the proof of Proposition 4.2) where balls are defined as

$$B_i = \{(\gamma, t) : |\gamma - \bar{\gamma}_i| + |t - \bar{t}_i| \leq R(N, L)\}$$

and again the norm  $|\gamma - \bar{\gamma}_i|$  is the Euclidian norm in  $\mathbb{R}^L$  while  $|t - \bar{t}_i|$  is the absolute value. Because  $B_d(0, 1) \times [-1, 1]$  is compact<sup>33</sup> it can be covered by a number of balls of the order of  $R(N, L)^{-(L+1)}$  (the extra dimension due to  $t$ ) and thus  $\bar{\mathfrak{N}}(N, L) \asymp R_N^{\max} d(\mathcal{B}_N)^L R(N, L)^{-(L+1)}$ .

The choice of  $R(N, L)$  is based on the same reasoning as before and the fact that the functions

$$(\gamma, t) \rightarrow \frac{\tilde{K}_{T_N}(s'\gamma - v) \cos(ty)}{\max\left(f_{\tilde{S}, \tilde{V}}(s, v), m_N\right)}$$

and

$$(\gamma, t) \rightarrow \mathbb{E} \left[ \frac{\tilde{K}_{T_N}(\tilde{S}'\gamma - \tilde{V}) \cos(tY)}{\max\left(f_{\tilde{S}, \tilde{V}}(\tilde{S}, \tilde{V}), m_N\right)} \right]$$

are Lipschitz with constant  $m_N^{-1} T_N^{L+2}$ . The same is true if we replace sin by cos. We can thus take the same  $t$ ,  $\omega$  and  $L$  and thus  $v_N$  as in the proof of Proposition 4.2. However due to the different covering number the constant  $C_1$  changes which yields a different constant  $M(\alpha)$ .

Q.E.D.

7.0.4. *Proof of Proposition 4.4.* First note that

$$\begin{aligned} (\widehat{f_\Delta} - f_\Delta)(\delta) &= - \int_{\mathbb{R}^L} f_{\Delta|\tilde{\Gamma}}(\delta|\gamma) f_{\tilde{\Gamma}}(\gamma) \mathbf{1}\{\gamma \in \mathcal{B}_N^c\} d\gamma \\ &\quad + \int_{\mathbb{R}^L} \left( \widehat{f_{\Delta|\tilde{\Gamma}}}(\delta|\gamma) - f_{\Delta|\tilde{\Gamma}}(\delta|\gamma) \right) f_{\tilde{\Gamma}}(\gamma) \mathbf{1}\{\gamma \in \mathcal{B}_N\} d\gamma \\ &\quad + \int_{\mathbb{R}^L} \widehat{f_{\Delta|\tilde{\Gamma}}}(\delta|\gamma) \left( \widehat{f_{\tilde{\Gamma}}}(\gamma) - f_{\tilde{\Gamma}}(\gamma) \right) \mathbf{1}\{\gamma \in \mathcal{B}_N\} d\gamma \end{aligned}$$

which yields

$$\begin{aligned} \left| (\widehat{f_\Delta} - f_\Delta)(\delta) \right| &\leq \int_{\mathbb{R}^L} f_{\Delta|\tilde{\Gamma}}(\delta|\gamma) f_{\tilde{\Gamma}}(\gamma) \mathbf{1}\{\gamma \in \mathcal{B}_N^c\} d\gamma \\ &\quad + \int_{\mathbb{R}^L} \left| \widehat{f_{\Delta|\tilde{\Gamma}}}(\delta|\gamma) - f_{\Delta|\tilde{\Gamma}}(\delta|\gamma) \right| f_{\tilde{\Gamma}}(\gamma) \mathbf{1}\{\gamma \in \mathcal{B}_N\} d\gamma \\ &\quad + \int_{\mathbb{R}^L} \widehat{f_{\Delta|\tilde{\Gamma}}}(\delta|\gamma) \left| \widehat{f_{\tilde{\Gamma}}}(\gamma) - f_{\tilde{\Gamma}}(\gamma) \right| \mathbf{1}\{\gamma \in \mathcal{B}_N\} d\gamma \\ &\leq \int_{\mathbb{R}^L} f_{\Delta|\tilde{\Gamma}}(\delta|\gamma) f_{\tilde{\Gamma}}(\gamma) \mathbf{1}\{\gamma \in \mathcal{B}_N^c\} d\gamma \\ &\quad + \int_{\mathbb{R}^L} \left| \widehat{f_{\Delta|\tilde{\Gamma}}}(\delta|\gamma) - f_{\Delta|\tilde{\Gamma}}(\delta|\gamma) \right| f_{\tilde{\Gamma}}(\gamma) \mathbf{1}\{\gamma \in \mathcal{B}_N\} d\gamma \\ &\quad + \left\| \frac{\widehat{f_{\tilde{\Gamma}}} - f_{\tilde{\Gamma}}}{f_{\tilde{\Gamma}}} \mathbf{1}\{\mathcal{B}_N\} \right\|_\infty \int_{\mathbb{R}^L} \left| \widehat{f_{\Delta|\tilde{\Gamma}}}(\delta|\gamma) \right| f_{\tilde{\Gamma}}(\gamma) \mathbf{1}\{\gamma \in \mathcal{B}_N\} d\gamma \end{aligned}$$

<sup>33</sup> $B_d(0, 1)$  is a Euclidian ball centered at 0 or radius 1.

$$\begin{aligned}
&\leq \int_{\mathbb{R}^L} f_{\Delta|\bar{\Gamma}}(\delta|\gamma) f_{\bar{\Gamma}}(\gamma) \mathbf{1}\{\gamma \in \mathcal{B}_N^c\} d\gamma \\
&\quad + \left(1 + \left\| \frac{\widehat{f_{\bar{\Gamma}}} - f_{\bar{\Gamma}}}{f_{\bar{\Gamma}}} \mathbf{1}\{\mathcal{B}_N\} \right\|_{\infty} \right) \int_{\mathbb{R}^L} \left| \widehat{f_{\Delta|\bar{\Gamma}}}(\delta|\gamma) - f_{\Delta|\bar{\Gamma}}(\delta|\gamma) \right| f_{\bar{\Gamma}}(\gamma) \mathbf{1}\{\gamma \in \mathcal{B}_N\} d\gamma \\
&\quad + \left\| \frac{\widehat{f_{\bar{\Gamma}}} - f_{\bar{\Gamma}}}{f_{\bar{\Gamma}}} \mathbf{1}\{\mathcal{B}_N\} \right\|_{\infty} \int_{\mathbb{R}^L} f_{\Delta|\bar{\Gamma}}(\delta|\gamma) f_{\bar{\Gamma}}(\gamma) \mathbf{1}\{\gamma \in \mathcal{B}_N\} d\gamma \\
&\leq \int_{\mathbb{R}^L} f_{\Delta|\bar{\Gamma}}(\delta|\gamma) f_{\bar{\Gamma}}(\gamma) \mathbf{1}\{\gamma \in \mathcal{B}_N^c\} d\gamma \\
&\quad + \left(1 + \left\| \frac{\widehat{f_{\bar{\Gamma}}} - f_{\bar{\Gamma}}}{f_{\bar{\Gamma}}} \mathbf{1}\{\mathcal{B}_N\} \right\|_{\infty} \right) \int_{\mathbb{R}^L} \left| \widehat{f_{\Delta|\bar{\Gamma}}}(\delta|\gamma) - f_{\Delta|\bar{\Gamma}}(\delta|\gamma) \right| f_{\bar{\Gamma}}(\gamma) \mathbf{1}\{\gamma \in \mathcal{B}_N\} d\gamma \\
&\quad + M_{\Delta} \left\| \frac{\widehat{f_{\bar{\Gamma}}} - f_{\bar{\Gamma}}}{f_{\bar{\Gamma}}} \mathbf{1}\{\mathcal{B}_N\} \right\|_{\infty}
\end{aligned}$$

thus, using the Cauchy-Schwartz inequality

$$\begin{aligned}
(\widehat{f_{\Delta}} - f_{\Delta})^2(\delta) &\leq 3 \int_{\mathbb{R}^L} f_{\Delta|\bar{\Gamma}}^2(\delta|\gamma) f_{\bar{\Gamma}}(\gamma) \mathbf{1}\{\gamma \in \mathcal{B}_N^c\} d\gamma \\
&\quad + 3 \left(1 + \left\| \frac{\widehat{f_{\bar{\Gamma}}} - f_{\bar{\Gamma}}}{f_{\bar{\Gamma}}} \mathbf{1}\{\mathcal{B}_N\} \right\|_{\infty} \right)^2 \int_{\mathbb{R}^L} \left( \widehat{f_{\Delta|\bar{\Gamma}}}(\delta|\gamma) - f_{\Delta|\bar{\Gamma}}(\delta|\gamma) \right)^2 f_{\bar{\Gamma}}(\gamma) \mathbf{1}\{\gamma \in \mathcal{B}_N\} d\gamma \\
&\quad + 3M_{\Delta}^2 \left\| \frac{\widehat{f_{\bar{\Gamma}}} - f_{\bar{\Gamma}}}{f_{\bar{\Gamma}}} \mathbf{1}\{\mathcal{B}_N\} \right\|_{\infty}^2 \\
&\leq 3M_{\Delta} \int_{\mathbb{R}^L} f_{\Delta|\bar{\Gamma}}(\delta, \gamma) \mathbf{1}\{\gamma \in \mathcal{B}_N^c\} d\gamma \\
&\quad + 3\|f_{\bar{\Gamma}}\|_{\infty} \left(1 + \left\| \frac{\widehat{f_{\bar{\Gamma}}} - f_{\bar{\Gamma}}}{f_{\bar{\Gamma}}} \mathbf{1}\{\mathcal{B}_N\} \right\|_{\infty} \right)^2 \int_{\mathbb{R}^L} \left( \widehat{f_{\Delta|\bar{\Gamma}}}(\delta|\gamma) - f_{\Delta|\bar{\Gamma}}(\delta|\gamma) \right)^2 \mathbf{1}\{\gamma \in \mathcal{B}_N\} d\gamma \\
&\quad + 3M_{\Delta}^2 \left\| \frac{\widehat{f_{\bar{\Gamma}}} - f_{\bar{\Gamma}}}{f_{\bar{\Gamma}}} \mathbf{1}\{\mathcal{B}_N\} \right\|_{\infty}^2.
\end{aligned}$$

The inequality is now obtained by integration over  $\delta$ .

Q.E.D.

7.0.5. *Proof of Proposition 4.5.* We introduce the notations

$$\begin{aligned}
\bar{f}_{\Delta|\bar{\Gamma}}(\delta|\gamma) &:= \frac{1}{2\pi} \int_{-\infty}^{\infty} K(t h_{N,\gamma}) e^{-i\delta t} \mathcal{F}_1 \left[ f_{\Delta|\bar{\Gamma}} \right] (t|\gamma) dt, \\
R(t, \gamma) &:= \frac{\mathbf{1} \left\{ \left| \mathcal{F}_1 \left[ \widehat{f_{Y_0, \bar{\Gamma}}} \right] (t, \gamma) \right| > r_{Y_0, N} \right\}}{\mathcal{F}_1 \left[ \widehat{f_{Y_0, \bar{\Gamma}}} \right] (t, \gamma)} - \frac{1}{\mathcal{F}_1 \left[ f_{Y_0, \bar{\Gamma}} \right] (t, \gamma)}.
\end{aligned}$$

The following decomposition holds at a fixed  $\gamma$  by means of the Plancherel identity

$$\begin{aligned} \left\| \left( \widehat{f_{\Delta|\tilde{\Gamma}}} - f_{\Delta|\tilde{\Gamma}} \right) (\cdot|\gamma) \right\|_2^2 &\leq 4 \left\| \left( \overline{f_{\Delta|\tilde{\Gamma}}} - f_{\Delta|\tilde{\Gamma}} \right) (\cdot|\gamma) \right\|_2^2 \\ &+ \frac{2}{\pi} \int_{-\infty}^{\infty} \frac{K(t h_{N,\gamma})^2}{\left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right|^2} \left| \left( \mathcal{F}_1 \left[ \widehat{f_{Y_0+\Delta, \tilde{\Gamma}}} \right] - \mathcal{F}_1 \left[ f_{Y_0+\Delta, \tilde{\Gamma}} \right] \right) (t, \gamma) \right|^2 dt \\ &+ \frac{2}{\pi} \int_{-\infty}^{\infty} K(t h_{N,\gamma})^2 |R(t, \gamma)|^2 \left| \left( \mathcal{F}_1 \left[ \widehat{f_{Y_0+\Delta, \tilde{\Gamma}}} \right] - \mathcal{F}_1 \left[ f_{Y_0+\Delta, \tilde{\Gamma}} \right] \right) (t, \gamma) \right|^2 dt \\ &+ \frac{2}{\pi} \int_{-\infty}^{\infty} K(t h_{N,\gamma})^2 \left| \mathcal{F}_1 \left[ f_{Y_0+\Delta, \tilde{\Gamma}} \right] (t, \gamma) \right|^2 |R(t, \gamma)|^2 dt. \end{aligned}$$

We conclude using Lemma 7.1 below and the fact that by conditional independence

$$\frac{\left| \mathcal{F}_1 \left[ f_{Y_0+\Delta, \tilde{\Gamma}} \right] (t, \gamma) \right|^2}{\left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right|^4} = \frac{\left| \mathcal{F}_1 \left[ f_{\Delta, \tilde{\Gamma}} \right] (t, \gamma) \right|^2}{\left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right|^2}.$$

Q.E.D.

Lemma 7.1 below is an adaptation of the lemma of Neumann (1997). Denote by

$$\psi(t, \gamma) := \frac{1}{\left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right|} \min \left( 1, \frac{r_{Y_0, N}}{\left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right|} \right).$$

**Lemma 7.1.**

$$\sup_{t \in [-R_N^{\max}, R_N^{\max}], \gamma \in \mathcal{B}_N} \{ \psi(t, \gamma)^{-1} |R(t, \gamma)| \} = O_p(1).$$

7.0.6. *Proof of Lemma 7.1.* We distinguish between two cases.

Case 1: Let  $t$  and  $\gamma$  be such that  $\left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right| < 2r_{Y_0, N}$ . Then,  $\psi(t, \gamma)^{-1} \leq 2 \left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right|$  and it suffices to upper bound in probability  $\left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right| |R(t, \gamma)|$ . By definition of  $R(t, \gamma)$ ,  $\left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right| |R(t, \gamma)| \leq 1$  on the event  $\left\{ \left| \mathcal{F}_1 \left[ \widehat{f_{Y_0, \tilde{\Gamma}}} \right] (t, \gamma) \right| \leq r_{Y_0, N} \right\}$ , while  $\left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right| |R(t, \gamma)| \leq (r_{Y_0, N})^{-1} \left| \left( \mathcal{F}_1 \left[ \widehat{f_{Y_0, \tilde{\Gamma}}} \right] - \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] \right) (t, \gamma) \right|$  on the complementary event  $\left\{ \left| \mathcal{F}_1 \left[ \widehat{f_{Y_0, \tilde{\Gamma}}} \right] (t, \gamma) \right| > r_{Y_0, N} \right\}$ . This yields

$$\sup_{(t, \gamma): \left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right| < 2r_{Y_0, N}} \{ \psi(t, \gamma)^{-1} |R(t, \gamma)| \} = O_p(1).$$

Case 2: Let now  $t$  and  $\gamma$  be such that  $\left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right| \geq 2r_{Y_0, N}$ . Then,  $\psi(t, \gamma)^{-1} \leq 2 (r_{Y_0, N})^{-1} \left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right|^2$  and it suffices to upper bound in probability  $(r_{Y_0, N})^{-1} \left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right|^2 |R(t, \gamma)|$ .

By definition of  $R(t, \gamma)$ ,

$$\begin{aligned} & (r_{Y_0, N})^{-1} \left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right|^2 |R(t, \gamma)| \\ & \leq (r_{Y_0, N})^{-1} \left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right| \left( \mathbf{1} \left\{ \left| \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} (t, \gamma) \right| \leq r_{Y_0, N} \right\} \right. \\ & \quad \left. + \frac{\left| \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} (t, \gamma) - \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right|}{\left| \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} (t, \gamma) \right|} \mathbf{1} \left\{ \left| \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} (t, \gamma) \right| > r_{Y_0, N} \right\} \right) \end{aligned}$$

Using

$$\frac{1}{\left| \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} (t, \gamma) \right|} \leq \frac{1}{\left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right|} + \frac{\left| \left( \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} - \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] \right) (t, \gamma) \right|}{\left| \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} (t, \gamma) \right| \left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right|},$$

we obtain

$$\begin{aligned} & (r_{Y_0, N})^{-1} \left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right|^2 |R(t, \gamma)| \\ & \leq (r_{Y_0, N})^{-1} \left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right| \left\{ \left| \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} (t, \gamma) \right| \leq r_{Y_0, N} \right\} \\ & + (r_{Y_0, N})^{-1} \left( \left| \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} (t, \gamma) - \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right| + \frac{\left| \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} (t, \gamma) - \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right|^2}{\left| \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} (t, \gamma) \right|} \right) \mathbf{1} \left\{ \left| \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} (t, \gamma) \right| > r_{Y_0, N} \right\} \\ & \leq (r_{Y_0, N})^{-1} \left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right| \left\{ \left| \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} (t, \gamma) \right| \leq r_{Y_0, N} \right\} \\ & + \left( (r_{Y_0, N})^{-1} \left| \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} (t, \gamma) - \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right| \right. \\ & \quad \left. + (r_{Y_0, N})^{-2} \left| \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} (t, \gamma) - \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right|^2 \right) \mathbf{1} \left\{ \left| \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} (t, \gamma) \right| > r_{Y_0, N} \right\}. \end{aligned}$$

From the definition of the upper bound on the rate  $r_{Y_0, N}$ , the last term in the sum is, uniformly in  $t$  and  $\gamma$  such that  $\left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right| \geq 2r_{Y_0, N}$ , bounded in probability.

Moreover, because  $\left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right| \geq 2r_{Y_0, N}$ ,

$$\begin{aligned} \mathbf{1} \left\{ \left| \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} (t, \gamma) \right| \leq r_{Y_0, N} \right\} & \leq \mathbf{1} \left\{ \left| \left( \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} - \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] \right) (t, \gamma) \right| \geq \left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right| - r_{Y_0, N} \right\} \\ & \leq \mathbf{1} \left\{ \left| \left( \widehat{\mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right]} - \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] \right) (t, \gamma) \right| \geq \left| \mathcal{F}_1 \left[ f_{Y_0, \tilde{\Gamma}} \right] (t, \gamma) \right| / 2 \right\} \end{aligned}$$

$$\leq 2 \frac{\left| \left( \mathcal{F}_1 \left[ \widehat{f_{Y_0, \bar{\Gamma}}} \right] - \mathcal{F}_1 \left[ f_{Y_0, \bar{\Gamma}} \right] \right) (t, \gamma) \right|}{\left| \mathcal{F}_1 \left[ f_{Y_0, \bar{\Gamma}} \right] (t, \gamma) \right|}$$

which yields

$$(r_{Y_0, N})^{-1} \left| \mathcal{F}_1 \left[ f_{Y_0, \bar{\Gamma}} \right] (t, \gamma) \right| \mathbf{1} \left\{ \left| \mathcal{F}_1 \left[ \widehat{f_{Y_0, \bar{\Gamma}}} \right] (t, \gamma) \right| \leq r_{Y_0, N} \right\} \leq (r_{Y_0, N})^{-1} \left| \left( \mathcal{F}_1 \left[ \widehat{f_{Y_0, \bar{\Gamma}}} \right] - \mathcal{F}_1 \left[ f_{Y_0, \bar{\Gamma}} \right] \right) (t, \gamma) \right|,$$

thus the first term is also, uniformly in  $t$  and  $\gamma$  such that  $\left| \mathcal{F}_1 \left[ f_{Y_0, \bar{\Gamma}} \right] (t, \gamma) \right| \geq 2r_{Y_0, N}$ , bounded in probability.

Q.E.D.

7.0.7. *Proof of Proposition 4.6.* The proposition follows from adapting the upper bounds in Comte and Lacour (2011), (4.17) and the assumptions made.

Q.E.D.

7.0.8. *Proof of Theorem 6.1.* Take  $\phi$  such that  $E[|\phi(Y_0)| + |\phi(Y_1)|] < \infty$  and  $s \in H^+$ , we have

$$\begin{aligned} \mathbb{E}[\zeta(D)\phi(Y)|S = s] &= \mathbb{E}[\zeta(D)\phi(Y)|S = s] \\ &= \mathbb{E}[\phi(Y_j)] - \mathbb{E}[\mathbf{1} \{s'\bar{\Gamma} > 0\} \phi(Y_j)] \quad (\text{using (A-4)}) \\ &= \mathbb{E}[\phi(Y_j)] - \int_{\mathbb{S}^L} \mathbf{1} \{s'\gamma > 0\} \overline{\mathbb{E}[\phi(Y_j)|\bar{\Gamma} = \cdot]}(\gamma) f_{\bar{\Gamma}}(\gamma) d\sigma(\gamma) \\ &= \mathbb{E}[\phi(Y_j)] - \mathcal{H} \left( \overline{\mathbb{E}[\phi(Y_j)|\bar{\Gamma} = \cdot]}(\gamma) f_{\bar{\Gamma}} \right) (s) \\ (7.14) \quad &= \frac{1}{2} \mathbb{E}[\phi(Y_j)] - \mathcal{H} \left( \overline{\mathbb{E}[\phi(Y_j)|\bar{\Gamma} = \cdot]} f_{\bar{\Gamma}} \right)^- (\gamma)(s) \end{aligned}$$

Let us give more details on how we obtain the last equality. We shall use classical results from harmonic analysis on the sphere that are recalled in Gautier and Kitamura (2009). A square integrable function on the sphere can be decomposed in a Fourier-Laplace series. The classical basis is a double indices sequence of functions  $(h_{nl})_{l=1, \dots, h(n, L), n=0, \dots, \infty}$  called the basis of spherical harmonics. It is composed of even and odd functions. Even functions are such that the index  $n$  is even and odd functions are such that  $n$  is odd.  $\left( \overline{\mathbb{E}[\phi(Y_j)|\bar{\Gamma} = \cdot]}(\gamma) f_{\bar{\Gamma}} \right)^-$  is the decomposition of  $\overline{\mathbb{E}[\phi(Y_j)|\bar{\Gamma} = \cdot]}(\gamma) f_{\bar{\Gamma}}$  on these odd basis functions. For every  $n = 2p$  for  $p \in \mathbb{N}$  and every  $l = 1, \dots, h(n, L)$ ,  $\int_{\mathbb{S}^L} h_{nl}(s) d\sigma(s) = 0$  and  $\mathcal{H}(h_{nl}) = 0$ . The function  $h_0$  ( $h(0, L) = 1$ ) is constant and equal to  $|\mathbb{S}^L|^{-1/2}$  (the value of the constant is such that the function is of norm 1). Take now a function  $f$ , by linearity and the results that we have recalled,

$$\mathcal{H}(f^+) = \int_{\mathbb{S}^L} \frac{f^+(\gamma)}{|\mathbb{S}^L|^{1/2}} d\sigma(\gamma) \mathcal{H} \left( \frac{1}{|\mathbb{S}^L|^{1/2}} \right)$$

$$\begin{aligned}
&= \int_{\mathbb{S}^L} \frac{f^+(\gamma)}{|\mathbb{S}^L|} d\sigma(\gamma) \mathcal{H}(1) \\
&= \int_{\mathbb{S}^L} \frac{f^+(\gamma)}{|\mathbb{S}^L|} d\sigma(\gamma) \frac{|\mathbb{S}^L|}{2} \quad (\text{because } \mathcal{H} \text{ is an integral over a hemisphere}) \\
&= \frac{1}{2} \int_{\mathbb{S}^L} f^+(\gamma) d\sigma(\gamma) \\
&= \frac{1}{2} \int_{\mathbb{S}^L} f(\gamma) d\sigma(\gamma) \quad (\text{because the odd part is orthogonal to } h_0).
\end{aligned}$$

Identity (7.14) yields that

$$(7.15) \quad \frac{1}{2} \mathbb{E}[\phi(Y_j)] - \mathbb{E}[\zeta(D)\phi(Y)|S = s] = \mathcal{H} \left( \overline{\mathbb{E}[\phi(Y_j)|\bar{\Gamma} = \cdot]}(\gamma) f_{\bar{\Gamma}} \right)^-(s).$$

The left hand-side in (7.15) is only defined in  $H^+$  while the right hand-side is defined on the whole sphere  $\mathbb{S}^L$  and is an odd function. Thus  $\frac{1}{2} \mathbb{E}[\phi(Y_j)] - \mathbb{E}[\zeta(D)\phi(Y)|S = s]$  can be extended in a natural way on the whole sphere as an odd function through

$$\begin{aligned}
\forall s \in H^+, \quad R_j(s) &= \frac{1}{2} \mathbb{E}[\phi(Y_j)] - \mathbb{E}[\zeta(D)\phi(Y)|S = s] \\
\forall s \in -H^+, \quad R_j(s) &= -R_j(-s).
\end{aligned}$$

Identity (6.2) follows by a simple manipulation of odd functions (see Gautier and Kitamura (2009)).

Identity (6.3) follows from the above discussion.

Let us now prove (6.4). Because of Assumption (A-3), the smoothing properties of the Hemispherical transform and the Sobolev embeddings, the right hand-side of (7.15) is continuous on the whole sphere. Therefore, the function  $R_j$ , which is only defined above on  $H^+ \cup (-H^+)$ , can be extended by continuity on  $\partial H^+$ . Because the extension should be an odd function, it satisfies, for any point  $\tilde{s}$  on  $\partial H^+$ ,  $R_j(\tilde{s}) = -R_j(-\tilde{s})$ . This yields (6.4).

Q.E.D.

## REFERENCES

- [1] AAKVIK, A., J. J. HECKMAN, AND E. J. VYTLACIL (2005): “Estimating Treatment Effects for Discrete Outcomes when Responses to Treatment Vary: an Application to Norwegian Vocational Rehabilitation Programs”. *Journal of Econometrics*, **125**, 15–51.
- [2] ABBRING, J. H., AND J. J. HECKMAN (2007): “Econometric Evaluation of Social Programs, Part III: Distributional Treatment Effects, Dynamic Treatment Effects, Dynamic Discrete Choice, and General Equilibrium Policy Evaluation”. *Handbook of Econometrics*, J.J. Heckman and E.E. Leamer (eds.), Vol. 6, North Holland, Chapter 72.
- [3] ABADIE, A., J. ANGRIST, AND G. IMBENS (2002): “Instrumental Variables Estimates of the Effect of Subsidized Training on the Quantiles of Trainee Earnings”. *Econometrica*, **70**, 91–117.



- [4] ANGRIST, J., K. GRADY AND G. IMBENS G. (2000): “The Interpretation of Instrumental Variables Estimators in Simultaneous Equations Models with an Application to the Demand for Fish”. *Review of Economic Studies*, **67**, 499–527.
- [5] BERAN, R., AND P. HALL (1992): “Estimating Coefficient Distributions in Random Coefficient Regressions”. *Annals of Statistics*, **20**, 1970–1984.
- [6] BERAN, R., A. FEUERVERGER, AND P. HALL (1996): “On Nonparametric Estimation of Intercept and Slope in Random Coefficients Regression”. *Annals of Statistics*, **24**, 2569–2592.
- [7] BJÖRKLUND, A., AND R. MOFFITT (1987): “The Estimation of Wage and Welfare Gains in Self-Selection Models”, *Review of Economics and Statistics*, **69**, 42–49.
- [8] BUTUCEA, C. (2004): “Deconvolution of Supersmooth Densities with Smooth Noise”. *Canadian Journal of Statistics*, **32**, 181–192.
- [9] BUTUCEA, C., AND A. B. TSYBAKOV (2007): “Sharp Optimality in Density Deconvolution with Dominating Bias. I”. *Rossiiskaya Akademiya Nauk. Teoriya Veroyatnostei i ee Primeneniya*, **52**, 111–128.
- [10] CARNEIRO, P., K. T. HANSEN, AND J. HECKMAN (2003): “Estimating Distributions of Treatment Effects With an Application to the Return to Schooling and Measurement of the Effect of Uncertainty on College Choice”. *International Economic Review*, **44**, 361–422.
- [11] CARRASCO, M., J. P. FLORENS, AND E. RENAULT (2007): “Linear Inverse Problems in Structural Econometrics Estimation Based on Spectral Decomposition and Regularization”. *Handbook of Econometrics*, J. J. Heckman and E. E. Leamer (eds.), vol. 6B, North Holland, chapter 77, 5633–5751.
- [12] CARRASCO, M., AND J. P. FLORENS (2011): “A Spectral Method for Deconvolving a Density”. *Econometric Theory*, **27**, 546–581.
- [13] CAVALIER, L. (2000): “Efficient Estimation of a Density in a Problem of Tomography”. *Annals of Statistics*, **28**, 630–647.
- [14] CAVALIER, L. (2001): “On the Problem of Local Adaptive Estimation in Tomography”. *Bernoulli*, **7**, 63–78.
- [15] CHERNOZHUKOV, V., AND C. HANSEN. (2005): “An IV Model of Quantile Treatment Effects”. *Econometrica*, **73**, 245–261.
- [16] COMTE, F., AND C. LACOUR (2011): “Data-driven Density Estimation in the Presence of Additive Noise with Unknown Distribution”. *Journal of the Royal Statistical Society. Series B*, **73**, 601–627.
- [17] DEVROYE, L. (1989): “Consistent Deconvolution in Density Estimation”. *Canadian Journal of Statistics*, **17**, 235–239.
- [18] ELBERS, C., AND G. RIDDER (1982): “True and Spurious Duration Dependence: The Identifiability of the Proportional Hazard Models”. *Review of Economics Studies*, **49**, 403–410.
- [19] EVANS, L. C. (1998): *Partial Differential Equations*. Graduate Studies in Mathematics, American Mathematical Society.
- [20] EVDOKIMOV, K. (2010): “Identification and Estimation of a Nonparametric Panel Data Model with Unobserved Heterogeneity”. Working paper.
- [21] EVDOKIMOV, K., AND H. WHITE (2011): “An Extension of a Lemma of Kotlarski”. Working paper.

- [22] FAN, Y., AND S. S. PARK (2010): “Sharp Bounds on the Distribution of Treatment Effects and Their Statistical Inference”. *Econometric Theory*, **26**, 931–951.
- [23] FAN, Y., AND D. ZHU (2009): “Partial Identification and Confidence Sets for Functional of the Joint Distribution of the Potential Outcomes”. Working paper.
- [24] FIRPO, S., AND G. RIDDER (2008): “Bounds on Functionals of the Distribution of Treatment Effects”. Working paper.
- [25] FRÉCHET, M. (1951): “Sur Les Tableaux de Corrélation Dont les Marges Sont Données”, *Annales de l'Université de Lyon, Série 3*, **14**, 53–77.
- [26] FOX, J., AND A. GANDHI (2011): “A Simple Nonparametric Approach to Estimating the Distribution of Random Coefficients in Structural Models”. Working Paper.
- [27] GAÏFFAS, S. (2005): “Convergence Rates for Pointwise Curve Estimation with a Degenerate Design”. *Mathematical Methods of Statistics*, **14**, 1-27.
- [28] GAÏFFAS, S. (2009): “Uniform Estimation of a Signal Based on Inhomogeneous Data”. *Statistica Sinica*, **19**, 427-447.
- [29] GAUTIER, E., AND Y. KITAMURA (2009): “Nonparametric Estimation in Random Coefficients Binary Choice Models”. Preprint [arXiv:0907.2451](https://arxiv.org/abs/0907.2451), forthcoming in *Econometrica*.
- [30] GAUTIER, E., AND E. LE PENNEC (2011): “Adaptive Estimation in Random Coefficients Binary Choice Models Using Needlet Thresholding”. Preprint [arXiv:1106.3503](https://arxiv.org/abs/1106.3503).
- [31] GUERRE, E. (1999): “Efficient Random Rates for Nonparametric Regression Under Arbitrary Designs”. Working Paper.
- [32] HALL, P., MARRON, J. S., NEUMANN, M. H., AND TETTERINGTON, D. M. (1997): “Curve Estimation When the Design Density is Low”. *Annals of Statistics*, **25**, 756-770.
- [33] HECKMAN, J. J., AND B. SINGER (1984): “A Method for Minimizing the Impact of Distributional Assumptions in Econometric Models for Duration Data”. *Econometrica*, **52**, 271–320.
- [34] HECKMAN, J. J., AND J. SMITH (1998): “Evaluating The Welfare State”. In: Strom, S. (Ed.), *Econometrics and Economic Theory of the Twentieth Century: The Ragnar Frisch Centennial Symposium*. Cambridge University Press, New York, pp. 241–318.
- [35] HECKMAN, J. J., J. SMITH, AND N. CLEMENTS (1997): “Making The Most Out Of Programme Evaluations and Social Experiments: Accounting For Heterogeneity in Programme Impacts”. *Review of Economic Studies*, **64**, 487–635.
- [36] HECKMAN, J. J., AND E. VYTLACIL (1999) “Local Instrumental Variables and Latent Variable Models for Identifying and Bounding Treatment Effects”. *Proceedings of the National Academy of Science, USA*, **96**, 4730–4734.
- [37] HECKMAN, J. J., AND E. VYTLACIL (2005): “Structural Equations, Treatment Effects, and Econometric Policy Evaluation”. *Econometrica*, **73**, 669–738.
- [38] HECKMAN, J. J., AND E. VYTLACIL (2007a) “Econometric Evaluation of Social Programs, Part I: Causal Models, Structural Models and Econometric Evaluation of Public Policies”. *Handbook of Econometrics*, J.J. Heckman and E.E. Leamer (eds.), Vol. 6, North Holland, Chapter 70.
- [39] HECKMAN, J. J., AND E. VYTLACIL (2007b) “Econometric Evaluation of Social Programs, Part II: Using the Marginal Treatment Effect to Organize Alternative Econometric Estimators to Evaluate Social Programs, and to

- Forecast their Effects in New Environments”. *Handbook of Econometrics*, J.J. Heckman and E.E. Leamer (eds.), Vol. 6, North Holland, Chapter 71.
- [40] HELGASON, S. (1999): *The Radon Transform*. 2nd edition. Birkhauser Boston.
- [41] HODERLEIN, S., J. KLEMELÄ, AND E. MAMMEN (2010): “Analyzing the Random Coefficient Model Nonparametrically”. *Econometric Theory*, **26**, 804–837.
- [42] Hoeffding, W. (1940): “Masstabinvariante Korrelationstheorie”. *Schriften des Mathematischen Instituts und Institutes Für Angewandte Mathematik der Universität Berlin*, **5**, 179–233.
- [43] HORVITZ, D. G., D. J. THOMPSON (1952): “A Generalization of Sampling Without Replacement From a Finite Universe”. *Journal of the American Statistical Association*, **47**, 663–685.
- [44] ICHIMURA, H., AND T. S. THOMPSON (1998): “Maximum Likelihood Estimation of a Binary Choice Model with Random Coefficients of Unknown Distribution”. *Journal of Econometrics*, **86**, 269–295.
- [45] IMBENS, G. W., AND J. D. ANGRIST (1994): “Identification and Estimation of Local Average Treatment Effects”. *Econometrica*, **62**, 467–475.
- [46] IMBENS, G. W., AND W. K. NEWEY (2009): “Identification and Estimation of Triangular Simultaneous Equations Models Without Additivity Corresponding”. *Econometrica*, **77**, 1481–1512.
- [47] JOHANNES, J. (2009): “Deconvolution with Unknown Error Distribution”. *Annals of Statistics*, **37**, 2301–2323.
- [48] KLEIN, T. (2010): “Heterogeneous Treatment Effects: Instrumental Variables Without Monotonicity?”. *Journal of Econometrics*, **155**, 99–116.
- [49] KOROSTELEV, A. P., AND A. B. TSYBAKOV (1993): *Minimax Theory of Image Reconstruction*. Springer, New-York, *Lecture Notes in Statistics* **82**.
- [50] LEWBEL, A. (2000): “Semiparametric Qualitative Response Model Estimation with Unknown Heteroscedasticity or Instrumental Variables”. *Journal of Econometrics*, **97**, 145–177.
- [51] LEWBEL, A. (2007): “Endogenous Selection or Treatment Model Estimation”. *Journal of Econometrics*, **141**, 777–806.
- [52] LOUNICI, K., AND R. NICKL (2011): “Global Uniform Risk Bounds for Wavelet Deconvolution Estimators”. *Annals of Statistics*, **39**, 201–231.
- [53] MAKAROV, G. D. (1981): “Estimates of the Distribution Function of a Sum of Two Random Variables when the Marginal Distributions are Fixed”. *Theory of Probability and its Applications*, **26**, 803–806.
- [54] MANSKI, C. F. (1997): “Monotone Treatment Effect”. *Econometrica*, **65**, 1311–1334.
- [55] MOFFITT, R. (2008): “Estimating Marginal Treatment Effects in Heterogeneous Populations”. *Annals of Economics and Statistics*, **91/92**, 239–261.
- [56] NATTERER, F. (1986): *The Mathematics of Computerized Tomography*. Wiley, Chichester.
- [57] NEUMANN, M. H. (1997): “On the Effect of Estimating the Error Density in Nonparametric Deconvolution”. *Journal of Nonparametric Statistics*, **7**, 307–330.
- [58] VAN DER VAART, A., AND J. WELLNER (1996): *Weak Convergence and Empirical Processes*. New York: Springer.
- [59] VYTLACIL, E. (2002): “Independence, Monotonicity, and Latent Index Models: An Equivalence Result”. *Econometrica*, **70** 331–341.

CREST (ENSAE), 3 AVENUE PIERRE LAROUSSE, 92 245 MALAKOFF CEDEX, FRANCE.

*E-mail address:* `Eric.Gautier@ensae-paristech.fr`

BOSTON COLLEGE, CHESTNUT HILL, MA 02467, USA.

*E-mail address:* `Stefan.Hoderlein@bc.edu`