

# Applied Causal Inference Powered by ML and AI

Victor Chernozhukov\*

Christian Hansen<sup>†</sup>

Nathan Kallus<sup>‡</sup>

Martin Spindler<sup>§</sup>

Vasilis Syrgkanis<sup>¶</sup>

February 28, 2024

Publisher: Online

Version 0.1.1

\* MIT

<sup>†</sup> Chicago Booth

<sup>‡</sup> Cornell University

<sup>§</sup> Hamburg University

<sup>¶</sup> Stanford University

# Unobserved Confounders, Instrumental Variables, and Proxy Controls

# 12

"Without Philip Wright  
would there have been causal DAGs?  
Who can really say?"

– Kei Hirano.\*

In this chapter we discuss various models with unobserved confounders, where the adjustment strategies we have discussed no longer work. We start with sensitivity analysis of causal inference to the presence of unobserved confounders. Then we discuss identification of causal effects when instrumental variables or proxy controls are available.

12.1 The Difficulty of Causal Inference with an Unobserved Confounder . . . . .	318
12.2 Impact of Confounders on Causal Effect Identification and Sensitivity Analysis . . .	319
12.3 Partially Linear IV Models . . . . .	322
A Wage Equation with Unobserved Ability . . . . .	322
Aggregate Market Demand . . . . .	324
Limits of Average Causal Effect Identification under Partial Linearity . . . . .	325
12.4 Nonlinear IV Models . . . . .	328
The LATE Model . . . . .	328
The IV Quantile Model* . . . . .	330
12.5 Partially Linear SEMs with Griliches-Chamberlain Proxy Controls . . . . .	331
12.6 Nonlinear Models with Proxy Controls* . . . . .	333
12.A Proofs . . . . .	337
Latent Confounder Bias Result: Theorem 12.2.1 . . . . .	337
Partially Linear Outcome IV Model: Theorem 12.3.2 . . . . .	338
Partially Linear Compliance IV Model: Theorem 12.3.3 . . . . .	338
Linear Proxy Model: Theorem 12.5.1 . . . . .	339

---

\* Sewall Wright, son, and Philip Wright, father, were responsible for some of the greatest ideas in causal inference. Sewall Wright invented causal path diagrams (linear DAGs), and Philip Wright wrote down DAGs for supply-demand equations, proposed IV methods for their identification, and even proposed weather conditions as instruments. Just one of these contributions would probably be enough to get a QJE publication in 1970s and later, but it was not good enough in 1926 or so. Philip Wright is a (causal) parent of Sewall Wright, so he is one of the causes of DAGs (hence the haiku).

## 12.1 The Difficulty of Causal Inference with an Unobserved Confounder

"All happy statisticians are happy in their own way; but all the unhappy ones are all alike — they all do causal inference with observational data". L. Tolstoy in Anna Karenina (Source: [Twitter](#))

Here we consider models with an unobserved confounding variable. The presence of unobserved confounding variables complicates identification of causal effects. Without further assumptions it is impossible to identify causal effects in a setting with unobserved confounding variables.

For example, consider the following two basic models shown in the margin figure, where we can think of  $Y$  as wages,  $D$  as education, and  $A$  as latent ability.

If  $A$  is not observed, the two models in Figures 12.1 and 12.2 are statistically indistinguishable from each other. In the first model  $D$  has a causal effect on  $Y$ , and in the second it does not. Even with strong restrictions, as in Gaussian linear SEMs, the observed correlation between  $D$  and  $Y$  can always be rationalized either as a causal effect of  $D$  on  $Y$  or the result of a common cause  $A$  (homework). This observation applies more generally. While we cannot precisely pin down causal effects in these cases, we can still learn about causal effects by performing sensitivity analysis if we are willing to assume a bound on the strength of unobserved confounders. We discuss a practical and intuitive approach to sensitivity analysis in Section 12.2.

We may also make progress in learning causal effects in the presence of unobserved confounders by considering the use of instrumental variables (IVs) – additional random vectors  $Z$  that create exogenous variation in  $D$ . This approach was introduced by Philip Wright in 1928 [1]. The use of instruments renders many linear ASEMs identifiable, allowing us to perform inference on structural effects  $D \rightarrow Y$ . Some nonlinear ASEMs also become identifiable, though identification still fails for completely unrestricted nonlinear models. We discuss the use of instruments in Sections 12.3-12.4.

A related set of problems is when we observe multiple proxy measurements of the latent confounder  $A$ . For example, we may observe  $S$ , the SAT score, and  $Q$ , the ACT score, which may both be proxies for latent confounder,  $A$ , ability. Note that conditioning on  $Q$  and  $S$  does not block the backdoor path  $Y \leftarrow A \rightarrow D$ . Hence we cannot use the regression adjustment

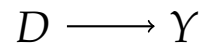


Figure 12.1:  $D$  causes  $Y$

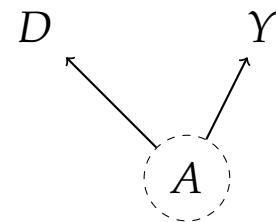


Figure 12.2:  $D$  and  $Y$  are caused by a latent factor  $A$

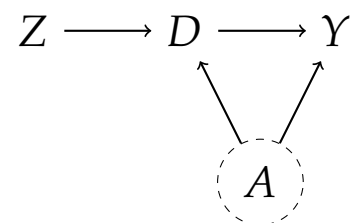


Figure 12.3: A DAG with Latent Confounder  $A$  and Instrument  $Z$ .

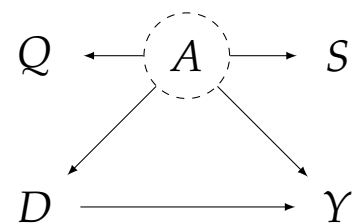


Figure 12.4: A DAG with two proxies for latent confounders.

method for identification of  $D \rightarrow Y$ . However, this problem is related to IVs, because we can effectively use one measurement in place of  $A$  and instrument it with another measurement to deal with the measurement error. This process can provide identification of the main effect  $D \rightarrow Y$ . In other words, we can use instrumental variable regression of  $Y$  on  $D$  and  $S$ , using  $D$  and  $Q$  as technical instrumental variables. This approach was introduced by Zvi Griliches in 1977 [2]. This model has also been extensively studied for nonlinear models as well, e.g., Miao et al. [3] and Deaner [4], especially in the recent literature. We discuss proxy approaches in Section 12.6.

## 12.2 Impact of Confounders on Causal Effect Identification and Sensitivity Analysis

**Example 12.2.1** (Partially Linear SEM) Consider the SEM

$$\begin{aligned} Y &:= \alpha D + \delta A + f_Y(X) + \epsilon_Y, \\ D &:= \gamma A + f_D(X) + \epsilon_D, \\ A &:= f_A(X) + \epsilon_A, \\ X &:= \epsilon_X, \end{aligned}$$

where, conditional on  $X$ ,  $\epsilon_Y, \epsilon_D, \epsilon_A$  are both mean zero and mutually uncorrelated. We further normalize

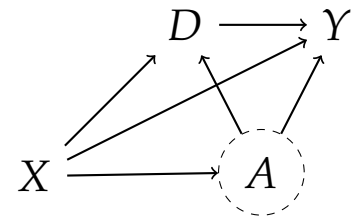
$$E[\epsilon_A^2] = 1.$$

The key structural parameter is  $\alpha$ :

$$\alpha = \partial_d Y(d)$$

where

$$Y(d) := (Y : do(D = d)).$$



**Figure 12.5:**  $X$  are observed confounders, and  $A$  are unobserved confounders.

To give context to our example, we can interpret  $Y$  as earnings,  $D$  as education,  $A$  as ability, and  $X$  as a set of observed background variables. In this example, we can interpret  $\alpha$  as the returns to schooling.

We start by applying the partialling out operator to get rid of the  $X$ 's in all of the equations. Define the partialling out operation of any random vector  $V$  with respect to another random vector  $X$  as the residual that is left after subtracting the best predictor

of  $V$  given  $X$ :

$$\tilde{V} = V - E[V | X].$$

If  $f$ 's are linear, we can replace  $E[V | X]$  by linear projection. After partialling out, we have a simplified system:

$$\begin{aligned} \tilde{Y} &:= \alpha \tilde{D} + \delta \tilde{A} + \epsilon_Y, \\ \tilde{D} &:= \gamma \tilde{A} + \epsilon_D, \\ \tilde{A} &:= \epsilon_A, \end{aligned}$$

where  $\epsilon_Y$ ,  $\epsilon_D$ , and  $\epsilon_A$  are uncorrelated.

Then the projection of  $\tilde{Y}$  on  $\tilde{D}$  recovers

$$\beta = E[\tilde{Y}\tilde{D}]/E[\tilde{D}^2] = \alpha + \phi,$$

where

$$\phi = \delta\gamma/E[(\gamma^2 + \epsilon_D^2)],$$

is the omitted confounder bias.

Omitted confounder bias is also often referred to as omitted variables bias.

The formula follows from inserting the expression for  $\tilde{D}$  into the definition of  $\beta$  and then simplifying the resulting expression using the assumptions on the  $\epsilon$ 's.

We can use this formula to bound  $\phi$  directly by making assumptions on the size of  $\delta$  and  $\gamma$ . An alternative approach can be based on the following characterization, based on partial  $R^2$ 's. This characterization essentially follows from Cinelli and Hazlett [5], with the slight difference that we have adapted the result to the partially linear model.<sup>1</sup>

1: [6] recently obtained a similar result for fully nonlinear models.

**Theorem 12.2.1** (Omitted Confounder Bias in Terms of Partial  $R^2$ 's) *In the setting given in Example 12.2.1,*

$$\phi^2 = \frac{R_{\tilde{Y}\sim\tilde{A}|\tilde{D}}^2 R_{\tilde{D}\sim\tilde{A}}^2}{(1 - R_{\tilde{D}\sim\tilde{A}}^2)} \frac{E[(\tilde{Y} - \beta\tilde{D})^2]}{E[(\tilde{D})^2]},$$

where  $R_{V\sim W|X}^2$  denotes the population  $R^2$  in the linear regression of  $V$  on  $W$ , after partialling out linearly  $X$  from  $V$  and  $W$ .

Therefore, if we place bounds on how much of the variation in  $\tilde{Y}$  and in  $\tilde{D}$  the unobserved confounder  $\tilde{A}$  is able to explain, we can bound the omitted confounder bias by

$$\sqrt{\phi^2}.$$

**Example 12.2.2** We consider an empirical example based on data surrounding the Darfur war. Specifically, we are interested in the effect of having experienced direct war violence on attitudes towards peace. The observed controls explain 12-15% of the variance of  $Y$ , beyond what's explained by the "treatment" variable, and 1% of the variance of treatment  $D$ . Therefore, suppose we are willing to accept that

$$R_{\tilde{Y} \sim \tilde{A} | \tilde{D}}^2 \leq .15, \quad R_{\tilde{D} \sim \tilde{A}}^2 \leq .01;$$

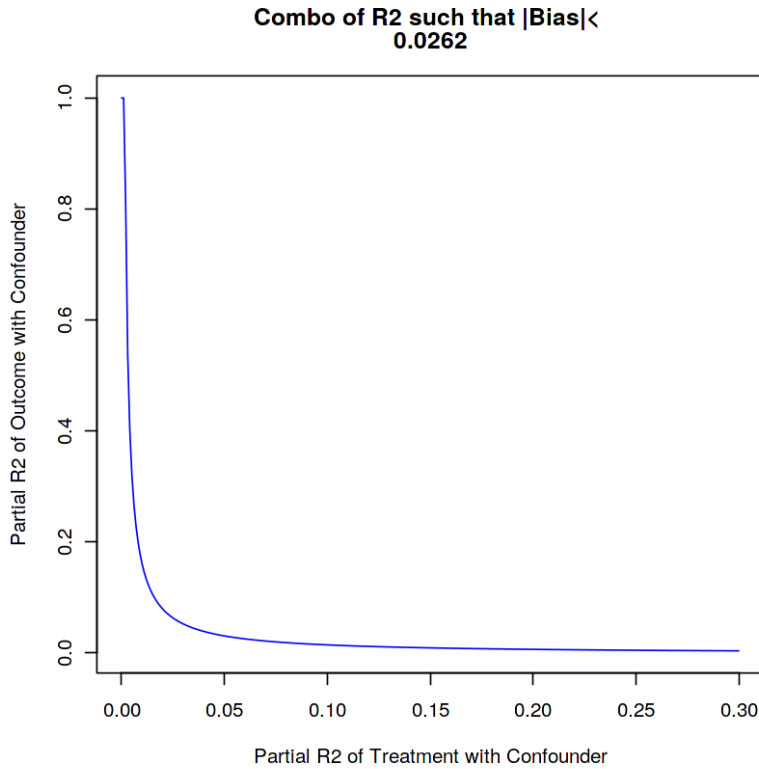
that is, we have a latent confounder that is no stronger than the observed controls for predicting  $Y$  and for predicting  $D$ .

Then, the upper/lower bound on  $\alpha$  is given by

$$\beta \pm \phi, \quad \phi^2 = \frac{.0015 \text{E}[(\tilde{Y} - \beta \tilde{D})^2]}{.99 \text{E}[\tilde{D}^2]}.$$

The estimated  $\beta$  is about .1. Plugging in estimates of  $\text{E}[(\tilde{Y} - \beta \tilde{D})^2]$  and  $\text{E}[(\tilde{D})^2]$  yields an estimated lower bound on  $\alpha$  of around .074. In Figure 12.6, we show the combination of all partial  $R^2$  such that the bias is less than .026. It shows that our conclusions about causal effects are not very sensitive to the presence of unknown confounders whose power is limited by the stated assumptions.

[DML Sensitivity R Notebook](#) carries out sensitivity analysis based on DML and the R package `Sensemkr` for the analysis of the Darfur wars data.



**Figure 12.6:** Sensitivity contour plots: The graph shows values of  $R^2_{\tilde{Y}|\tilde{D}|\tilde{A}}$  and  $R^2_{\tilde{D}|\tilde{A}}$  that give a given value of the bias  $|\hat{\phi}| = .026$ .

### 12.3 Partially Linear IV Models

When instrumental variables are available, it becomes possible to point identify causal effects in partially linear models and certain types of causal effects in nonlinear models. Here we begin with partially linear models.

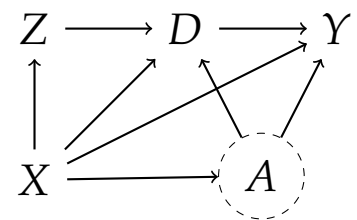
#### A Wage Equation with Unobserved Ability

**Example 12.3.1** (Returns to Education with Omitted Ability; Generalization of Griliches, 1977 [2]) Consider the ASEM

$$\begin{aligned} Y &:= \alpha D + \delta A + f_Y(X) + \epsilon_Y, \\ D &:= \beta Z + \gamma A + f_D(X) + \epsilon_D, \\ Z &:= f_Z(X) + \epsilon_Z, \\ A &:= f_A(X) + \epsilon_A, \\ X &:= \epsilon_X, \end{aligned}$$

where, conditional on  $X$ ,  $\epsilon_Y, \epsilon_D, \epsilon_Z, \epsilon_A$  have mean zero and are mutually uncorrelated.

We can interpret  $Y$  as earnings,  $D$  as education,  $A$  as ability,  $Z$



**Figure 12.7:** An IV model with observed and unobserved confounders.

as an observed shifter of education, and  $X$  as a set of observed background variables. The key structural parameter is  $\alpha$ , the returns to schooling, i.e.

$$\alpha = \partial_d Y(d),$$

where

$$Y(d) = Y : do(D = d).$$

Examples of instruments for schooling,  $Z$ , that have appeared in the literature include

- ▶ distance to college (Card [7]),
- ▶ compulsory schooling laws (Angrist [8]),
- ▶ offer to participate/offer to treat in a training program (many studies), and
- ▶ subsidies to finance education (Griliches, Heckman).

We apply the partialling-out operator to get rid of the  $X$ 's in all of the equations. As before, we define the partialling out operation of any random vector  $V$  with respect to another random vector  $X$  as the residual that is left after subtracting the best predictor of  $V$  given  $X$ :

$$\tilde{V} = V - E[V | X].$$

If  $f$ 's are linear, we replace  $E[V | X]$  with linear projection.

After partialling-out, we have a simplified system.

$$\begin{aligned} \tilde{Y} &:= \alpha \tilde{D} + \delta \tilde{A} + \epsilon_Y, \\ \tilde{D} &:= \beta \tilde{Z} + \gamma \tilde{A} + \epsilon_D, \\ \tilde{Z} &:= \epsilon_Z, \\ \tilde{A} &:= \epsilon_A, \end{aligned}$$

where  $\epsilon_Y, \epsilon_D, \epsilon_Z$ , and  $\epsilon_A$  are uncorrelated.

We immediately obtain the following result:

**Theorem 12.3.1** *In Example 12.3.1, we can rewrite an econometric measurement model for identification of  $\alpha$ :*

$$\tilde{Y} := \alpha \tilde{D} + U, \quad U \perp \tilde{Z},$$

where  $U = \delta \tilde{A} + \epsilon_Y$ . Alternatively, we can equivalently identify  $\alpha$  using the moment restriction

$$E[(\tilde{Y} - \alpha \tilde{D})\tilde{Z}] = 0.$$



The identification of  $\alpha$  follows from solving this equation,

$$\alpha = E[\tilde{Y}\tilde{Z}]/E[\tilde{D}\tilde{Z}],$$

provided the instruments are relevant:  $E[\tilde{D}\tilde{Z}] \neq 0$  or  $\beta \neq 0$ .

**Remark 12.3.1** (Neyman Orthogonality and DML) The target parameter  $\alpha$  is Neyman orthogonal with respect to nuisance parameters – the regression functions  $E[Y | X]$ ,  $E[D | X]$ , and  $E[Z | X]$ . Therefore we can use DML for learning and performing statistical inference on the parameter  $\alpha$ .

### Wright’s Causal Path Derivation

Starting from the DAG given in Figure 12.7, we obtain Figure 12.8 after partialling out.

Philip Wright (1928) [1] observed that the structural parameter  $\beta\alpha$ , the effect  $\tilde{Z} \rightarrow \tilde{Y}$ , is identified from the projection of  $\tilde{Y} \sim \tilde{Z}$ :

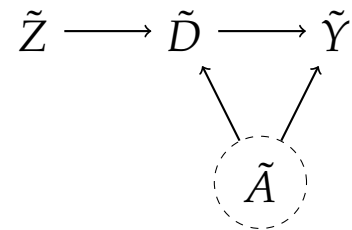
$$\beta\alpha = E[\tilde{Y}\tilde{Z}]/E[\tilde{Z}^2].$$

The structural parameter  $\beta$ , the effect of  $Z \rightarrow D$ , is identified from the projection of  $\tilde{D} \sim \tilde{Z}$ :

$$\beta = E[\tilde{D}\tilde{Z}]/E[\tilde{Z}^2].$$

$\alpha$ , the effect of  $D \rightarrow Y$ , is then identified by the ratio of the two provided  $\beta \neq 0$ :

$$\alpha = \frac{\beta\alpha}{\beta} = E[\tilde{Y}\tilde{Z}]/E[\tilde{D}\tilde{Z}].$$

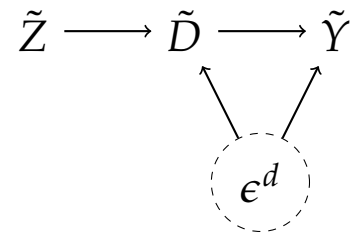


**Figure 12.8:** DAG corresponding to Figure 12.7 after partialling out observed confounder  $X$ .

### Aggregate Market Demand

Let’s apply our approach to a canonical example in economics: the identification of the price elasticity of demand using a supply shifter as an instrument.

**Example 12.3.2** (Market Demand; Generalization of P. Wright,



**Figure 12.9:** A DAG for aggregate demand, with the latent node  $\epsilon^d$  representing the demand shock

1928 [1]) Consider the ASEM

$$\begin{aligned} Y &:= \alpha D + f_Y(X) + \epsilon^d, \\ D &:= \beta Z + f_D(X) + \rho \epsilon^d + \gamma \epsilon^s, \\ Z &:= f_Z(X) + \epsilon_Z \end{aligned}$$

where  $\epsilon^d$ ,  $\epsilon^s$  and  $\epsilon_Z$  are mean zero and uncorrelated conditional on  $X$ . In this example,  $Y$  is (log) demand,  $D$  is (log) price,  $Z$  is an observed supply shifter,  $X$  is a vector of observed demand shifters,  $\epsilon^d$  is a demand shock, and  $\epsilon^s$  is a supply shock. The key parameter is  $\alpha$ , the price elasticity of demand:

$$\alpha = \partial_d Y(d),$$

where  $Y(d) := (Y : do(D = d))$ . Here we focus on only the demand side of the market and do not attempt to explicitly model the supply side.

Example 12.3.2 is equivalent to the previous Example 12.3.1 – set  $A = \epsilon^d$ ,  $\epsilon_Y = 0$ ,  $\epsilon^s = \epsilon_D$ , and so on. Hence, the identification method is the same as before.

In econometrics, the set-up here is sometimes referred to as a *limited information* model or formulation because we are focusing on identifying only a single equation in a more complicated underlying system.

## Limits of Average Causal Effect Identification under Partial Linearity

The result in Theorem 12.3.1 extends beyond the partially linear setting presented in Example 12.3.1 to the following non-linear structural equation model:

**Example 12.3.3** (Partially Linear Outcome IV Model) Consider the ASEM

$$\begin{aligned} Y &:= g_Y(\epsilon_Y)D + f_Y(A, X, \epsilon_Y), \\ D &:= f_D(Z, X, A, \epsilon_D), \\ Z &:= f_Z(X, \epsilon_Z), \\ A &:= f_A(X, \epsilon_A), \\ X &:= \epsilon_X, \end{aligned}$$

where,  $\epsilon_Y, \epsilon_D, \epsilon_Z, \epsilon_A$  are exogenous and mutually independent. The key structural parameter is:

$$\alpha := E[\partial_d Y(d)] = E[g_Y(\epsilon_Y)],$$

where

$$Y(d) = Y : do(D = d).$$

This parameter is typically referred to as the average marginal effect of the treatment.

Theorem 12.3.1 extends almost as is to this more general non-linear structural equation model.

**Theorem 12.3.2** *In Example 12.3.3, we can identify  $\alpha$  using the moment restriction*

$$E[(\tilde{Y} - \alpha \tilde{D})\tilde{Z}] = 0.$$

*The identification of  $\alpha$  follows from solving this equation,*

$$\alpha = E[\tilde{Y}\tilde{Z}] / E[\tilde{D}\tilde{Z}],$$

*provided the instruments are relevant:  $E[\tilde{D}\tilde{Z}] \neq 0$ .*

Note that the non-linear structural equation model in Example 12.3.3 imposes extra assumptions on the structural response function of the outcome  $Y$ . Thus our identification argument imposes more conditions on the structural equations than the ones that can be encoded via a DAG. Such auxiliary assumptions are required for identification of average treatment effects with instruments.

In particular, the identification argument relies on the fact that the unobserved confounder  $A$  enters in an additively separable manner in the outcome equation. If for instance,  $A$  was an input to the function  $g$ , i.e.  $Y := g_Y(A, \epsilon_Y)D + f_Y(A, X, \epsilon_Y)$ , then the quantity identified by the moment restriction in Theorem 12.3.2 would not correspond to an average treatment effect. In this case, the unobserved confounder creates heterogeneity in the treatment effect and also heterogeneity in the effect of the instrument on the treatment, typically referred to as the "compliance" (i.e., the correlation between  $Z$  and  $D$  varies with  $A$ ). This property is what renders the ratio quantity  $\alpha = E[\tilde{Y}\tilde{Z}] / E[\tilde{D}\tilde{Z}]$ , invalid for the causal estimand of interest.

In fact, it is the joint heterogeneity in both the outcome relationship and the compliance relationship that causes the problem. We show next that we could allow for a much more complex outcome model as long as the effect of the instrument on the treatment (compliance) is not heterogeneous in  $A$  or  $X$ .

**Example 12.3.4** (Partially Linear Compliance IV Model) Con-

sider the ASEM

$$\begin{aligned} Y &:= g_Y(A, X, \epsilon_Y)D + f_Y(A, X, \epsilon_Y), \\ D &:= g_D(\epsilon_D)Z + f_D(X, A, \epsilon_D), \\ Z &:= f_Z(X) + \epsilon_Z, \\ A &:= f_A(X, \epsilon_A), \\ X &:= \epsilon_X, \end{aligned}$$

where,  $\epsilon_Y, \epsilon_D, \epsilon_Z, \epsilon_A$  are exogenous and mutually independent. The key structural parameter is:

$$\alpha := E[\partial_d Y(d)] = E[g(A, X, \epsilon_Y)],$$

where

$$Y(d) = Y : do(D = d).$$

**Theorem 12.3.3** *In Example 12.3.4, we can identify  $\alpha$  using the moment restriction*

$$E[(\tilde{Y} - \alpha \tilde{D})\tilde{Z}] = 0.$$

*The identification of  $\alpha$  follows from solving this equation,*

$$\alpha = E[\tilde{Y}\tilde{Z}] / E[\tilde{D}\tilde{Z}],$$

*provided the instruments are relevant:  $E[\tilde{D}\tilde{Z}] \neq 0$ .*

Thus, we see that we need that either the effect of education on wages is not heterogeneous in the unobserved ability variable  $A$  or that the effect of the observed education shifter  $Z$  (e.g. distance to college) on education  $D$  is not heterogeneous in the unobserved ability variable to use the identification strategies presented in this section in the context of our education example. In Section 12.4, we will investigate what causal quantities are identifiable even in non-linear structural equation models, where the unobserved confounder creates heterogeneity in both the treatment effect and in the compliance behavior.

**Remark 12.3.2** (Effect heterogeneity based on observables)

We note that allowing for  $X$  to enter the  $g_Y$  or  $g_D$  function in Example 12.3.3 and Example 12.3.4 (i.e. allowing for the treatment effect or compliance, i.e. effect of instrument on treatment, to vary with  $X$ ), is a more benign extension because  $X$  is an observed variable. In this case, we can repeat the identification strategies in this section, conditional on  $X$ , and

we can show with similar arguments that

$$\beta(X) := E[\partial_d Y(d) | X] = \frac{E[\tilde{Y}\tilde{Z} | X]}{E[\tilde{D}\tilde{Z} | X]}. \quad (12.3.1)$$

Then we can simply average these conditional estimates to get the average marginal effect:

$$\alpha = E[\beta(X)]. \quad (12.3.2)$$

Such an identification strategy was initiated in [9, 10] and was also recently used in the context of DML estimators [11–13]. In particular, the following moment condition that identifies  $\alpha$ ,

$$E\left[\beta(X) + \frac{(\tilde{Y} - \beta(X)\tilde{D})\tilde{Z}}{E[\tilde{D}\tilde{Z} | X]} - \alpha\right] = 0, \quad (12.3.3)$$

is Neyman orthogonal with respect to the nuisance functions  $\beta(X)$  and  $\gamma(X) := E[\tilde{D}\tilde{Z} | X]$ . We note that this identification strategy remains valid even if in Example 12.3.4 the instrument equation is fully non-linear, i.e.  $Z := f_Z(X, \epsilon_Z)$ .

## 12.4 Nonlinear IV Models

Once we consider nonlinear models, identification becomes a much more delicate matter. We first consider the local average treatment effect (LATE) model, and then we turn to quantile models.

### The LATE Model

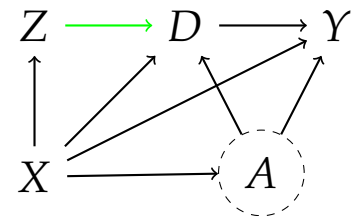
An important nonlinear IV model is the local average treatment effect model (LATE), proposed by Imbens and Angrist [14].

**Example 12.4.1 (LATE)** Consider the SEM, where

$$\begin{aligned} Y &:= f_Y(D, X, A, \epsilon_Y) \\ D &:= f_D(Z, X, A, \epsilon_D) \in \{0, 1\}, \\ Z &:= f_Z(X, \epsilon_Z) \in \{0, 1\}, \\ X &:= \epsilon_X, \quad A = \epsilon_A, \end{aligned}$$

where  $\epsilon$ 's are all independent, and

$z \mapsto f_D(z, A, X, \epsilon_D)$  is weakly increasing (weakly monotone).



**Figure 12.10:** LATE models. Green arrow denotes a monotone functional relation.

Suppose the instrument  $Z$  is an offer to participate in a training program and that  $D$  is the actual endogenous participation in the training program. Participation in the program may depend on unobservables  $A$ , such as ability or perseverance, that also affect the eventual outcome  $Y$ . We can also have background exogenous covariates  $X$  in the model.

Define

$$Y(d) := f_Y(d, X, A, \epsilon_Y) \text{ and } D(z) := f_D(z, X, A, \epsilon_D)$$

as the potential outcomes that result from applying fix-interventions in the corresponding equations from Example 12.4.1.

The model allows us to identify the local average treatment effect (LATE), defined as

$$\theta = E[Y(1) - Y(0) \mid D(1) > D(0)],$$

where  $\{D(1) > D(0)\}$  is the compliance event, where switching instrument value from  $Z = 0$  to  $Z = 1$  induces participation. Therefore LATE measures the average treatment effect conditional on compliance.

**Theorem 12.4.1** *In the LATE model, we have that  $\theta$  is identified by the ratio of two statistical parameters,*

$$\theta = \theta_1 / \theta_2,$$

where

$$\theta_1 := E[E[Y \mid X, Z = 1] - E[Y \mid X, Z = 0]],$$

and

$$\theta_2 := E[E[D \mid X, Z = 1] - E[D \mid X, Z = 0]],$$

provided that the instrument  $Z$  is relevant,  $\theta_2 > 0$ , and  $Z$  has full conditional support – namely  $0 < P(Z = 1 \mid X) < 1$ . Moreover,  $\theta_2$  identifies the probability of compliance:

$$\theta_2 = P[D(1) > D(0)].$$

The result has an intuitive interpretation.<sup>2</sup> In the event of compliance, the instrument moves the treatment as if experimentally, which induces quasi-experimental variation in the outcome. We measure the probability of compliance with  $\theta_2$

2: In the model with no  $X$  the ratio  $\theta_1 / \theta_2$  is equivalent to Wright's [1] IV estimand.

and the average induced changes in outcome by  $\theta_1$ . Taking the ratio is then like conditioning on the compliance event. See the proof in Section 12.A for details.

The ratio can be recognized as the ratio of average treatment effects of  $Z$  on  $Y$  and  $D$ ,

$$\theta_1 = ATE(Z \rightarrow Y),$$

$$\theta_2 = ATE(Z \rightarrow D).$$

This assertion follows from the application of the backdoor criterion. Therefore in order to perform inference on LATE, we can simply re-use the tools for performing inference on two ATEs.

**Remark 12.4.1** (DML for  $\theta_1/\theta_2$ ) We can apply DML to obtain  $\hat{\theta}_1$  and  $\hat{\theta}_2$  and then construct the estimator  $\hat{\theta} = \hat{\theta}_1/\hat{\theta}_2$  via the plug-in principle. This approach automatically has the Neyman orthogonality property.

### The IV Quantile Model\*

Another nonlinear IV model is the following model that exploits monotonicity in the unobservable shock in the outcome equation to obtain identification.

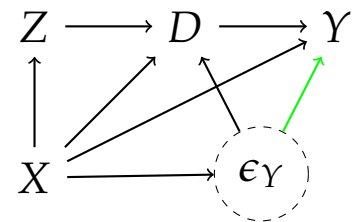
**Example 12.4.2** (IV Quantile Model) Consider the SEM

$$\begin{aligned} Y &= f_Y(D, X, \epsilon_Y), \\ D &= f_D(Z, X, \epsilon_Y, \epsilon_D), \\ Z &= f_Z(X, \epsilon_Z), \\ X &= \epsilon_X, \end{aligned}$$

where  $\epsilon$ 's are all independent,

$$f_Y(D, X, \cdot) : [0, 1] \mapsto \mathbb{R} \text{ is strictly increasing,}$$

and  $\epsilon_Y$  is normalized to have uniform distribution on  $(0, 1)$ . The context could be given from the demand example, where  $Y$  is demand,  $D$  price,  $\epsilon_Y$  a demand shock,  $\epsilon_D$  a supply shock;  $X$  the set of background variables, and  $Z$  a set of instrumental variables. The function  $f_Y(d, x, u)$  is the  $u$ -th quantile of the structural function of  $f_Y(d, x, \epsilon_Y)$ , which is the demand function in this context. For example,  $f_Y(d, x, 1/2)$  is the median structural function.



**Figure 12.11:** IV Quantile Model. The green arrow represents a strictly monotonic effect.

The testable implication of the IV Quantile Model is the following.

**Theorem 12.4.2** *In the IV Quantile Model, the testable moment restriction is*

$$P[Y \leq f_Y(D, X, u) \mid Z, X] = u,$$

for each  $u \in (0, 1)$ . There exist regularity conditions, analogous to instrument relevance, under which the structural function  $f_Y$  is identified from this restriction.

In practice, linear forms  $f_Y(D, X, u) = \alpha(u)'D + \beta(u)'X$  are often used. Adopting a linear functional form leads to method of moments approaches such as the IV quantile regression for performing inference on structural quantile functions.

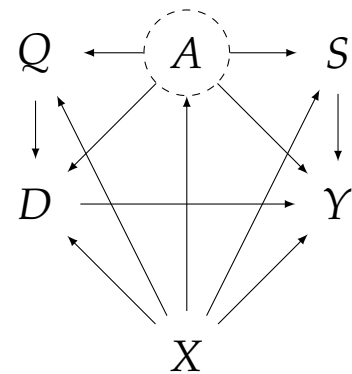
**Remark 12.4.2** (DML for IVQR Models) The problem of constructing DML for IVQR problems is considered open. Neyman-orthogonal approaches for the partially linear IVQR models are sketched out in the review [15] and may be a good place to start.

Code for IV Quantile Models can be found [here](#).

## 12.5 Partially Linear SEMs with Griliches-Chamberlain Proxy Controls

Suppose we are interested in the causal effect of college education on earnings in the presence of an unobserved confounder – individual ability. Here we show that we can recover the effect of college education on earnings in the presence of latent ability using proxies for ability, but not the effect of ability itself.

**Example 12.5.1** (Earnings with Omitted Ability; Griliches, 1977 [2]; Griliches and Chamberlain, 1977 [16]) Consider the



**Figure 12.12:** A DAG with Controls and Proxy Controls



ASEM

$$\begin{aligned}
 Y &:= \alpha D + \delta A + \iota S + f_Y(X) + \epsilon_Y, \\
 D &:= \gamma A + \beta Q + f_D(X) + \epsilon_D, \\
 Q &:= \eta A + f_Q(X) + \epsilon_Q, \\
 S &:= \phi A + f_S(X) + \epsilon_S, \\
 A &:= f_A(X) + \epsilon_A, \\
 X &:= \epsilon_X,
 \end{aligned}$$

where  $\epsilon_Y, \epsilon_D, \epsilon_Q, \epsilon_S, \epsilon_A, \epsilon_X$  have mean zero and are uncorrelated conditional on  $X$ . Interpret  $Y$  as earnings,  $D$  as college degree,  $A$  as ability,  $Q$  and  $S$  as proxies of ability, and  $X$  as a set of observed background variables. Example proxies  $Q$  and  $S$  are

- $Q$  is test scores or grades in some period  $t_0$  and  $S$  is test scores or grades at a later period  $t_1$ .

The key structural parameter is  $\alpha$ , the returns to schooling; i.e.

$$\alpha = \partial_d Y(d),$$

where  $Y(d) = Y : d_0(D = d)$ .

After partialling out we are left with the DAG in Figure 12.13:

$$\begin{aligned}
 \tilde{Y} &:= \alpha \tilde{D} + \delta \tilde{A} + \iota \tilde{S} + \epsilon_Y, \\
 \tilde{D} &:= \gamma \tilde{A} + \beta \tilde{Q} + \epsilon_D, \\
 \tilde{Q} &:= \eta \tilde{A} + \epsilon_Q, \\
 \tilde{S} &:= \phi \tilde{A} + \epsilon_S, \\
 \tilde{A} &:= \epsilon_A,
 \end{aligned}$$

where  $\epsilon_Y, \epsilon_D, \epsilon_Q, \epsilon_S, \epsilon_A$  are uncorrelated. The idea now is to replace  $\tilde{A}$  in the equation for  $\tilde{Y}$  with  $\tilde{S}$ . Note that because  $S$  enters the  $Y$  equation directly, we cannot consider using  $\tilde{Q}$  to proxy for  $\tilde{A}$ . We still cannot learn  $\alpha$  from the regression of  $\tilde{Y}$  on  $\tilde{D}$  and  $\tilde{S}$  though as  $S$  is an imperfect proxy for  $A$ . The following result, which provides an IV approach to identify  $\alpha$ , is immediate via substitution.<sup>3</sup>

**Theorem 12.5.1** *Assume that all variables in Example 12.5.1 are square-integrable. Then we have the following measurement equation:*

$$\begin{aligned}
 \tilde{Y} &= \alpha \tilde{D} + \bar{\delta} \tilde{S} + U, \quad E[U(\tilde{D}, \tilde{Q})] = 0, \\
 U &= -\delta \epsilon_S / \phi + \epsilon_Y; \quad \bar{\delta} = \iota + \delta / \phi.
 \end{aligned}$$

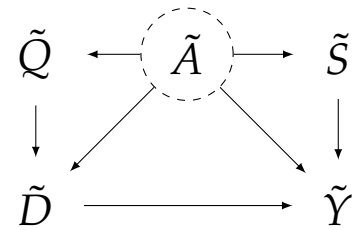


Figure 12.13: A DAG with Proxy Controls After Partialling Out

3: Prove the result as a reading exercise. Substitute  $\tilde{A} = (\tilde{S} - \epsilon_S) / \phi$  in the first equation and use the assumptions on the disturbances.

Here  $\alpha$  is identified from the moment condition  $E[U(\tilde{D}, \tilde{Q})] = 0$ , which is equivalent to using  $\tilde{Q}$  as an instrument for  $\tilde{S}$ , provided that  $\tilde{D}$  and the best linear predictor of  $\tilde{S}$  using  $\tilde{Q}$  and  $\tilde{D}$  have non-degenerate covariance matrix.

Note that  $\tilde{Q}$  here plays the role of a *technical* instrument for  $\tilde{S}$ . This approach recovers  $\alpha$ , but not  $\delta$ . For inference, we can employ the DML method for IV models; see also Chapter 13.

**Remark 12.5.1** (Neyman Orthogonality and DML) The formulation of the target parameter given above is Neyman-orthogonal, and high-quality estimation and statistical inference can be carried out using DML. In essence, we just residualize the system, using cross-fitted residuals, and then apply standard instrumental variable methods from econometrics to perform inference on the structural parameter of interest.

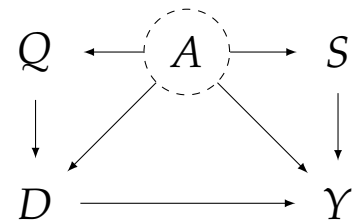
## 12.6 Nonlinear Models with Proxy Controls\*

An important recent development is "proximal causal inference," which generalizes early work by Griliches and Chamberlain [16].<sup>†</sup>

**Example 12.6.1** (Miao, Geng, and Tchetgen Tchetgen [3]) We consider the following model encoded in the DAG in Figure 12.14:

$$\begin{aligned} Y &:= f_Y(D, S, A, \epsilon_Y), \\ D &:= f_D(A, Q, \epsilon_D), \\ Q &:= f_Q(A, \epsilon_Q), \\ S &:= f_S(A, \epsilon_S), \\ A &:= \epsilon_A, \end{aligned}$$

where  $\epsilon$ 's are mutually independent. We can endow the same context to this model as in Example 12.5.1.



**Figure 12.14:** A SEM with Proxy Controls  $Q$  and  $S$ . Note that conditioning on  $Q$  and  $S$  does not block the backdoor path  $Y \leftarrow A \rightarrow D$ , hence we cannot use the regression adjustment method for identification of  $D \rightarrow Y$ .

<sup>†</sup> The most relevant papers include, amongst others, the stream of work by Tchetgen Tchetgen and collaborators, as well as the dissertation work of Deaner. Here we describe some results of the first group specialized to the discrete case.

Here we can introduce background exogenous controls  $X$  in each of the equations, but we don't do so to save notation. Notice that the model in Example 12.6.1 generalizes the Example 12.5.1 to the nonparametric case.

**Assumption 12.6.1** *In Example 12.6.1, assume*

- (a) *Variables  $Q$ ,  $S$  and  $A$  are finitely discrete and take on the same number of values.*
- (b) *The matrix  $\Pi(S | Q, d)$ , whose  $s^{\text{th}}$  row and  $q^{\text{th}}$  column is  $p(s | q, d)$ , is invertible for each value  $d$ .*

Condition (b) is analogous to the usual relevance condition in IV and basically says that the two proxies  $S$  and  $Q$  have sufficient joint variation at any value of  $d$  to allow  $Q$  to serve as an "instrument" for  $S$ . The discreteness assumption can be generalized to a more general completeness condition; see Miao et al.[3] and Deaner [4]. As with the usual IV relevance condition, Condition (b) is testable from the data. In contrast, the DAG itself and the other conditions involve an unobserved variable  $A$  and are therefore generally untestable. The validity of these untestable conditions must be assessed using contextual knowledge about the empirical problem.

**Theorem 12.6.1** *Under Assumption 12.6.1,  $p(y : do(d))$  is identifiable by the proximal formula:*

$$p(y : do(d)) = \Pi(y | d, Q) \Pi(S | Q, d)^{-1} \Pi(S), \quad (12.6.1)$$

*where  $\Pi(y | d, Q)$  and  $\Pi(S)$  are row and column vectors whose entries are of the form  $p(y | d, q)$  and  $p(s)$ .*

The mnemonic way to think about the formula above is that we are doing a kind of instrumental regression of  $Y$  on  $S$ , while instrumenting  $S$  with  $Q$ , which is exactly how we dealt with the linear version of this problem in Section 12.6.

**Remark 12.6.1** [17] and [18] provide moment functions defined in terms of efficient influence functions, which possess the Neyman orthogonality property, for estimating of the average treatment effect within this proxy control setting in the presence of a high-dimensional set of control variables. These moment functions can thus serve as the foundation for the use of DML inference methods for the average treatment effect in such settings.

## Notebooks

- ▶ [DML Sensitivity R Notebook](#) analyses the sensitivity of the DML estimate in the Darfur wars example to unobserved confounders using the Sensemakr package in R. [DML Sensitivity Python Notebook](#) does the same analysis in Python.
- ▶ [DML for Partially Linear IV R Notebook](#) and [DML for Partially Linear IV Python Notebook](#) carry out the DML IV analysis of the Acemoglu-Johnson-Robinson example, which considers the impact of the quality of institutions on economic growth, instrumenting quality of institutions with settler mortality. The notebook explores the partially linear IV model and tests for the presence of weak instruments. See Chapter 14 for further discussion of this example as well as discussion of weak identification/instruments.
- ▶ [DML for LATE Models R Notebook](#) and [DML for LATE Models Python Notebook](#) estimate the Local Average Treatment Effects of 401(K) participation on net financial wealth.
- ▶ [DML for Linear Proxy Controls Python Notebook](#) provides an application of using proxy controls to estimate the effect of smoking on birth weight.

## Study Problems

1. Explain omitted confounder bias to a fellow student (one paragraph). Explore using sensitivity analysis to aid in understanding robustness of economic conclusions to the presence of unobserved confounders in an empirical example of your choice. The [DML Sensitivity R Notebook](#) can be a helpful starting point but apply the ideas to a different empirical example. (You could use any of the previous examples we have analyzed).
2. Write a brief explanation of the idea of the instrumental variables regression model that would be appropriate for educating a fellow student. Discuss the idea of identifying the causal effect in this setting via path analysis in the spirit of what Philip Wright did. Illustrate your discussion with an empirical example. For example, revisit the analysis in [DML for Partially Linear IV R Notebook](#).

3. (Simulation.) Create a notebook to simulate one of the linear IV or proxy controls models that we've described. Assume there are no  $X$ 's for simplicity. Demonstrate numerically why using least squares may not be appropriate due to unobserved confounding. Demonstrate numerically how using instrumental variable regression overcomes the issue.
4. (LATE etc.) Explain to a fellow student in writing one of the nonlinear models (e.g. LATE, IV quantile model, or the nonlinear model with proxy controls) and how causal parameters in these models are identified. [DML for LATE Models R Notebook](#) could be a starting point for explaining LATE and illustrating your explanation with empirical results. (If you have a good empirical example for proxy controls, please let us know.)

## 12.A Proofs

### Latent Confounder Bias Result: Theorem 12.2.1

The proof heavily relies on the Frisch-Waugh-Lovell partialling out theorem (FWL) and the normalization on the variance of the latent confounder:

$$E[\tilde{A}^2] = 1. \quad (12.A.1)$$

The proof also relies on the properties of  $R_{U \sim V}^2$  which measures the proportion of variance of centered random variable  $U$  that is linearly explained by another centered random variable  $V$ :

$$R_{U \sim V}^2 = \frac{E\beta^2 V^2}{E[U^2]} = 1 - \frac{E[\epsilon^2]}{E[U^2]} = \frac{(E[UV])^2}{E[U^2]E[V^2]} = \text{Cor}^2(U, V),$$

where  $\beta = E[VU]/E[V^2]$  is the coefficient of the best linear projection of  $U$  onto  $V$ ,  $\epsilon = U - \beta V$  is the projection residual, and  $\text{Cor}(U, V)$  denotes the correlation between  $U$  and  $V$ . Note that  $R^2$  is symmetric in  $U$  and  $V$ :  $R_{U \sim V}^2 = R_{V \sim U}^2$ .

By FWL and the normalization (12.A.1), we have

$$\gamma = E[\tilde{A}\tilde{D}], \quad \delta = E[\tilde{A}\tilde{Y}]/E[\tilde{A}^2],$$

where

$$\begin{aligned} \tilde{Y} &= \tilde{Y} - \beta\tilde{D}; & \beta &= E[\tilde{Y}\tilde{D}]/E[\tilde{D}^2]; \\ \tilde{A} &= \tilde{A} - \tilde{\beta}\tilde{D}; & \tilde{\beta} &= E[\tilde{A}\tilde{D}]/E[\tilde{D}^2]. \end{aligned}$$

It follows that

$$\phi^2 = \frac{\gamma^2 \delta^2}{(E[\tilde{D}^2])^2} = \frac{(E[\tilde{A}\tilde{D}])^2 (E[\tilde{Y}\tilde{A}])^2}{(E[\tilde{D}^2])^2 (E[\tilde{A}^2])^2}.$$

Then the result follows from the normalization (12.A.1) and the following relations:

$$(E[\tilde{D}\tilde{A}])^2 = \text{Cor}^2(\tilde{D}, \tilde{A})E[\tilde{D}^2] = R_{\tilde{D} \sim \tilde{A}}^2 E[\tilde{D}^2],$$

$$(E[\tilde{Y}\tilde{A}])^2 = \text{Cor}^2(\tilde{Y}, \tilde{A})E[\tilde{Y}^2]E[\tilde{A}^2] = R_{\tilde{Y} \sim \tilde{A}}^2 E[\tilde{Y}^2]E[\tilde{A}^2],$$

$$E[\tilde{A}^2] = 1 - R_{\tilde{A} \sim \tilde{D}}^2 = 1 - R_{\tilde{D} \sim \tilde{A}}^2.$$

and noting that by definition  $R_{\tilde{Y} \sim \tilde{A}}^2 = R_{\tilde{Y} \sim \tilde{A}|\tilde{D}}^2$ .

□

**Partially Linear Outcome IV Model:****Theorem 12.3.2**

First note that since  $E[\tilde{Z} | X] = 0$ , we can re-write the moment condition as

$$E[(Y - \alpha D)\tilde{Z}] = 0.$$

We can use the structural equation for  $Y$  to replace  $Y$  in the moment equation:

$$E[(g_Y(\epsilon_Y)D + f_Y(A, X, \epsilon_Y) - \alpha D)\tilde{Z}] = 0.$$

Furthermore, since  $\tilde{Z} \perp\!\!\!\perp A, \epsilon_Y | X$ , we have that

$$\begin{aligned} E[f_Y(A, X, \epsilon_Y)\tilde{Z}] &= E[f_Y(A, X, \epsilon_Y)E[\tilde{Z} | X, A, \epsilon_Y]] \\ &= E[f_Y(A, X, \epsilon_Y)E[\tilde{Z} | X]] = 0. \end{aligned}$$

Thus we can re-write the moment equation as

$$E[(g_Y(\epsilon_Y)D - \alpha D)\tilde{Z}] = 0.$$

Solving for  $\alpha$  and using the fact that  $\epsilon_Y \perp\!\!\!\perp \tilde{Z}$ , we get

$$\alpha = \frac{E[g_Y(\epsilon_Y)D\tilde{Z}]}{E[D\tilde{Z}]} = \frac{E[g_Y(\epsilon_Y)]E[D\tilde{Z}]}{E[D\tilde{Z}]} = E[g_Y(\epsilon_Y)].$$

□

**Partially Linear Compliance IV Model:****Theorem 12.3.3**

Using the exact same arguments as in the proof of Theorem 12.3.2, we can deduce that the solution to the moment restriction takes the form

$$\alpha = \frac{E[g_Y(X, A, \epsilon_Y)D\tilde{Z}]}{E[D\tilde{Z}]} = \frac{E[g_Y(X, A, \epsilon_Y)E[D\tilde{Z} | X, A, \epsilon_Y]]}{E[D\tilde{Z}]}.$$

We now use the assumptions on the structural response functions of  $D$  and  $Z$  to argue that  $E[D\tilde{Z} | X, A, \epsilon_Y] = E[D\tilde{Z}]$ , i.e. the covariance of  $D$  and  $Z$  (aka compliance) is independent of  $X, A, \epsilon_Y$ . This independence would then imply the theorem, since we would get that

$$\alpha = E[g_Y(X, A, \epsilon_Y)].$$

First, we use the assumption on the structural response function of  $D$ :

$$E[D\tilde{Z} \mid X, A, \epsilon_Y] = E[(g_D(\epsilon_D)Z + f_D(X, A, \epsilon_D))\tilde{Z} \mid X, A, \epsilon_Y].$$

Using the fact that  $Z \perp\!\!\!\perp A, \epsilon_D, \epsilon_Y \mid X$ , and that  $E[\tilde{Z} \mid X] = 0$ , we can remove the term  $f_D(X, A, \epsilon_D)$  from the above equation:

$$E[D\tilde{Z} \mid X, A, \epsilon_Y] = E[g_D(\epsilon_D)Z\tilde{Z} \mid X, A, \epsilon_Y].$$

Using the additively separable assumption on the structural response of  $Z$  and the fact that  $\epsilon_Z$  is an exogenous independent variable, we have

$$\begin{aligned} E[D\tilde{Z} \mid X, A, \epsilon_Y] &= E[g_D(\epsilon_D)Z\epsilon_Z \mid X, A, \epsilon_Y] \\ &= E[g_D(\epsilon_D)\epsilon_Z^2 \mid X, A, \epsilon_Y] = E[g_D(\epsilon_D)\epsilon_Z^2] \end{aligned}$$

where we used the fact that all noise variables  $\epsilon_D, \epsilon_Y, \epsilon_Z$  are exogenous and mutually independent.  $\square$

### Linear Proxy Model: Theorem 12.5.1.

We substitute  $\tilde{A} = (\tilde{S} - \epsilon_S)/\phi$  in the equation  $\tilde{Y} := \alpha\tilde{D} + \delta\tilde{A} + \iota\tilde{S} + \epsilon_Y$  to obtain

$$\begin{aligned} \tilde{Y} &= \alpha\tilde{D} + \bar{\delta}\tilde{S} + U, \\ U &= -\delta\epsilon_S/\phi + \epsilon_Y; \quad \bar{\delta} = \iota + \delta/\phi. \end{aligned}$$

To verify

$$E[U] = 0$$

we observe using repeated substitutions that:

- ▶  $\tilde{D}$  is a linear combination of  $(\epsilon_A, \epsilon_Q, \epsilon_D)$ ,
- ▶  $\tilde{Q}$  is a linear combination of  $\epsilon_A$  and  $\epsilon_Q$ .
- ▶  $U$  is a linear combination of  $(\epsilon_S, \epsilon_Y)$ .

The conclusion follows from the assumption that

$$(\epsilon_A, \epsilon_Q, \epsilon_D, \epsilon_S, \epsilon_Y)$$

are all uncorrelated. The conclusion that  $\alpha$  is identified provided that  $\tilde{D}$  and the best linear predictor of  $\tilde{S}$  using  $\tilde{Q}$  and  $\tilde{D}$  have non-degenerate covariance matrices is left as an exercise.

$\square$



**The LATE Result: Theorem 12.4.1**

We can use, for example, the backdoor criterion to conclude that

$$E[E[D | Z = z, X]] = E[E[D(z) | X]] = ED(z).$$

Similarly,

$$E[E[Y | Z = z, X]] = E[E[Y(D(z)) | X]] = E[Y(D(z))].$$

Furthermore, by monotonicity, we have both

$$\theta_2 = E[D(1) - D(0)] = P(D(1) > D(0))$$

and

$$\begin{aligned} \theta_1 &= E[Y(D(1)) - Y(D(0))] \\ &= E[\{Y(1) - Y(0)\}1\{D(1) > D(0)\}]. \end{aligned}$$

Therefore

$$\theta_1/\theta_2 = E[Y(1) - Y(0) | D(1) > D(0)].$$

□

**Testable Restriction for the IV Quantile Model:  
Theorem 12.4.2**

The result is immediate from (i) the equivalence of the event  $Y \leq f_Y(D, X, u)$  and the event  $\epsilon_Y \leq u$ , which holds under the strict monotonicity assumption, and (ii) the independence of  $\epsilon_Y$  from  $Z$  and  $X$  which follows from the stated independence conditions. Using (i) and (ii), we have

$$\begin{aligned} P[Y \leq f_Y(D, X, u) | Z, X] &= P[\epsilon_Y \leq u | Z, X] \\ &= P[\epsilon_Y \leq u] = P[U(0, 1) \leq u] = u. \end{aligned}$$

□

**Identification in the Nonlinear Proxy Variables  
Model: Theorem 12.6.1**

To sketch a proof, the DAG implies that the observed variables  $D, Y, Q, S$  and the unobserved variable  $A$  obey the two

conditional independence relations:

$$(i) S \perp\!\!\!\perp (Q, D) \mid A \quad (ii) Q \perp\!\!\!\perp Y \mid (A, D). \quad (12.A.2)$$

These in turn imply

$$\begin{aligned} \Pi(S \mid Q, d) &= \Pi(S \mid Q, A, d)\Pi(A \mid Q, d) \\ &= \Pi(S \mid A)\Pi(A \mid Q, d) \end{aligned}$$

and

$$\begin{aligned} \Pi(y \mid Q, d) &= \Pi(y \mid Q, A, d)\Pi(A \mid Q, d) \\ &= \Pi(y \mid A, d)\Pi(A \mid Q, d). \end{aligned}$$

We now want to solve these equations for  $\Pi(y \mid A, d)$  in terms of quantities that could be learned in the data.

We will need invertibility of  $\Pi(S \mid Q, d)$  which requires invertibility of both  $\Pi(S \mid A)$  and  $\Pi(A \mid Q, d)$ . Under these invertibility conditions, we have

$$\Pi(A \mid Q, d) = \Pi(S \mid A)^{-1}\Pi(S \mid Q, d)$$

and

$$\Pi(y \mid Q, d) = \Pi(y \mid A, d)\Pi(S \mid A)^{-1}\Pi(S \mid Q, d),$$

which yield

$$\Pi(y \mid A, d) = \Pi(y \mid Q, d)\Pi(S \mid Q, d)^{-1}\Pi(S \mid A).$$

Next, because  $A$  blocks backdoor paths between  $D$  and  $Y$ , we have that

$$p(y \mid a : do(d)) = p(y \mid a, d) \quad (12.A.3)$$

or, after integrating out  $a$ ,

$$p(y : do(d)) = \Pi(y \mid A, d)\Pi(A),$$

which can be further expressed as

$$\Pi(y \mid d, Q) \Pi(S \mid Q, d)^{-1} \Pi(S), \quad (12.A.4)$$

using the derivations above.  $\square$

# Bibliography

- [1] Philip G. Wright. *The Tariff on Animal and Vegetable Oils*. New York: The Macmillan company, 1928 (cited on pages 318, 324, 325, 329).
- [2] Zvi Griliches. 'Estimating the returns to schooling: Some econometric problems'. In: *Econometrica: Journal of the Econometric Society* 45.1 (1977), pp. 1–22 (cited on pages 319, 322, 331).
- [3] Wang Miao, Zhi Geng, and Eric J. Tchetgen Tchetgen. 'Identifying causal effects with proxy variables of an unmeasured confounder'. In: *Biometrika* 105.4 (2018), pp. 987–993 (cited on pages 319, 333, 334).
- [4] Ben Deaner. 'Proxy controls and panel data'. In: *arXiv preprint arXiv:1810.00283* (2018) (cited on pages 319, 334).
- [5] Carlos Cinelli and Chad Hazlett. 'Making sense of sensitivity: Extending omitted variable bias'. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 82.1 (2020), pp. 39–67 (cited on page 320).
- [6] Victor Chernozhukov, Carlos Cinelli, Whitney Newey, Amit Sharma, and Vasilis Syrgkanis. 'Omitted Variable Bias in Machine Learned Causal Models'. In: *arXiv preprint arXiv:2112.13398* (2021) (cited on page 320).
- [7] David Card. 'Using Geographic Variation in College Proximity to Estimate the Return to Schooling'. In: *Aspects of Labor Market Behavior: Essays in Honor of John Vanderkamp*. Ed. by L. N. Christofides and R. Swidinsky. Toronto: University of Toronto Press, 1995, pp. 201–222 (cited on page 323).
- [8] Joshua D. Angrist and Alan B. Krueger. 'Does Compulsory School Attendance Affect Schooling and Earnings?'. In: *Quarterly Journal of Economics* 106.4 (1991), pp. 979–1014 (cited on page 323).
- [9] Alberto Abadie. 'Semiparametric instrumental variable estimation of treatment response models'. In: *Journal of econometrics* 113.2 (2003), pp. 231–263 (cited on page 328).
- [10] Peter M. Aronow and Allison Carnegie. 'Beyond LATE: Estimation of the Average Treatment Effect with an Instrumental Variable'. In: *Political Analysis* 21.4 (2013), pp. 492–506. (Visited on 02/27/2023) (cited on page 328).

- [11] Ryo Okui, Dylan S. Small, Zhiqiang Tan, and James M. Robins. 'DOUBLY ROBUST INSTRUMENTAL VARIABLE REGRESSION'. In: *Statistica Sinica* 22.1 (2012), pp. 173–205. (Visited on 02/27/2023) (cited on page 328).
- [12] Susan Athey and Stefan Wager. 'Policy learning with observational data'. In: *Econometrica* 89.1 (2021), pp. 133–161 (cited on page 328).
- [13] Vasilis Syrgkanis, Victor Lei, Miruna Oprescu, Maggie Hei, Keith Battocchi, and Greg Lewis. 'Machine learning estimation of heterogeneous treatment effects with instruments'. In: *Advances in Neural Information Processing Systems* 32 (2019) (cited on page 328).
- [14] Guido W. Imbens and Joshua D. Angrist. 'Identification and Estimation of Local Average Treatment Effects'. In: *Econometrica* 62.2 (1994), pp. 467–475 (cited on page 328).
- [15] Victor Chernozhukov, Christian Hansen, and Kaspar Wuthrich. 'Instrumental variable quantile regression'. In: *arXiv preprint arXiv:2009.00436* (2020) (cited on page 331).
- [16] Gary Chamberlain and Zvi Griliches. *More on brothers*. In *"Kinometrics: Determinants of Socioeconomic Success Within and Between Families"* (P. Taubman, Ed.) 1977 (cited on pages 331, 333).
- [17] Yifan Cui, Hongming Pu, Xu Shi, Wang Miao, and Eric Tchetgen Tchetgen. 'Semiparametric proximal causal inference'. In: *arXiv preprint arXiv:2011.08411* (2020) (cited on page 334).
- [18] Nathan Kallus, Xiaojie Mao, and Masatoshi Uehara. 'Causal inference under unmeasured confounding with negative controls: A minimax learning approach'. In: *arXiv preprint arXiv:2103.14029* (2021) (cited on page 334).