

Econometric Causality

The Central Role of Thought Experiments

James Heckman and Rodrigo Pinto*

Revised December 1, 2023

Abstract

This paper examines the econometric causal model and the interpretation of empirical evidence based on thought experiments that was developed by Ragnar Frisch and Trygve Haavelmo. We compare the econometric causal model with two currently popular causal frameworks: the Neyman-Rubin causal model and the Do-Calculus. The Neyman-Rubin causal model is based on the language of potential outcomes and was largely developed by statisticians. Instead of being based on thought experiments, it takes statistical experiments as its foundation. The Do-Calculus, developed by Judea Pearl and co-authors, relies on Directed Acyclic Graphs (DAGs) and is a popular causal framework in computer science and applied mathematics. We make the case that economists who uncritically use these frameworks often discard the substantial benefits of the econometric causal model to the detriment of more informative analyses. We illustrate the versatility and capabilities of the econometric framework using causal models developed in economics.

*James Heckman is the Henry Schultz Distinguished Service Professor of Economics and Public Policy at the University of Chicago; and Director of the Center for Economics of Human Development. Rodrigo Pinto is an Assistant Professor in the Department of Economics at the University of California – Los Angeles. We have benefitted from comments by Ed Vytlačil and three anonymous referees. This research was supported in part by NIH grant NICHD R37HD065072 and a grant from an anonymous funder. The views expressed in this paper are solely those of the authors and do not necessarily represent those of the funders or the official views of the National Institutes of Health.

Key words: Structural Equation Models, Causality, Causal Inference, Directed
Acyclic Graphs, Simultaneous Causality

JEL codes: C10, C18

Rodrigo Pinto
University of California at Los Angeles
Department of Economics
315 Portola Plaza, Room 8385
Los Angeles, CA 90095
(310) 825-0849
rodrig@econ.ucla.edu

James Heckman
The University of Chicago
Department of Economics
1126 E. 59th St.
Chicago, IL 60637
(773) 702-0634
jjh@uchicago.edu

1 Introduction

Sound economic and policy analysis is causal analysis. It analyzes the factors that produce outcomes and the role of various factors and policies in doing so. It quantifies policy impacts. It elucidates the mechanisms producing outcomes in order to understand how they operate, how they can be transported to different environments, how programs might be improved and which, if any, alternative mechanisms might be used to generate desired outcomes. It organizes evidence in interpretable frameworks. It uses all available information to give good policy advice and explicitly recognizes any limitations of data or models.

Good economic science systematically explores possible counterfactual worlds. It is grounded in thought experiments – what might happen if determinants of outcomes are changed. Credible hypotheticals are developed, analyzed, and tested with real-world data.

Models and thought experiments are central to economic analysis. Persons trained in economic theory or in the natural sciences routinely use them. Statisticians and computer scientists have only recently come to grips with the causal questions that have long been investigated by econometric pioneers such as Ragnar Frisch and Trygve Haavelmo. As a result, private languages and procedures designed to do part of what rigorous econometric models do have been developed without manifesting understanding of the full corpus of econometric theory, often refusing to cite it and reinventing portions of it.

These private languages bear the marks of their recent birth: concepts are often not precisely defined, and the conceptually distinct issues of definition of counterfactuals, their identification, and their estimation are often tangled together. In some fields heavily influenced by statistics, certain estimation techniques are claimed to be central for the definition or identification of counterfactuals when, in fact, they are only devices for recovering counterfactuals from data.

The current state of affairs would be of little concern if applied economists continued to draw on and extend the standard econometric model. Sadly, this is not the case. Many

econometricians and applied economists now emulate what they read in statistics and computer science journals. They have forgotten or never learned their own field’s foundational work to the detriment of rigorous causal analysis and testing among alternatives.

This paper discusses econometric causal analysis and recently developed causal models in fields outside economics. Our goal is to enhance the theory and practice of economic policy analysis by testing and synthesizing evidence, as well as interpreting it. This involves acquainting economists with a rich econometric legacy and situating recently advocated causal frameworks within the broader context of the econometric model.¹

The topic is broad and our paper is necessarily brief. We discuss some main points and illustrate them with analyses of a few prototypical economic models used to address policy problems and interpret evidence. It is impossible to convey here all of the insights of rigorous econometrics developed in the past 90 years.

This paper unfolds in the following way. We first define the notion of causality within a model. The concept is simple, but requires thought processes outside of statistics that are, nonetheless, quite intuitive. We discuss four distinct classes of policy problems that are addressed in econometric causal analyses. Some of them are either ignored or only partly addressed in the recent non-economic causal literatures. We demonstrate the conceptual clarity of the econometric approach and contrast it with that of rival approaches.

In particular, we consider two causal frameworks often advocated by statisticians and computer scientists. The first is the Neyman-Rubin model (1923; 1958; 1974; 1986; 1996), “NR” henceforward. It uses notions developed in rigorous econometrics but goes only part way toward implementing the full set of tools in the econometric approach to causal analysis and the interpretation of empirical evidence. It has important limitations for posing or analyzing routine policy problems outside a narrow “treatment-control” paradigm. It ignores

¹In this paper we focus on policy analysis but our message applies to a broader class of problems. The models developed in this paper also apply to the tasks of hypothesis testing, statistical inference, and synthesis of empirical evidence into interpretable evidence. Formulating meaningful alternatives is central to power analysis or Bayesian tests among alternatives.

the simultaneous equations model - a major achievement of econometrics - and replaces it with a litany of “confounding biases” readily addressed in rigorous econometrics. We also consider an approach to counterfactuals developed in computer science (“*do-calculus*,” [Pearl, 2012](#)), henceforth “DoC,” that relies critically on directed acyclic graphs (DAGs) and statistical conditional independence relationships. We demonstrate its limited capacity to address many important economic questions and address important empirical problems or to utilize many standard econometric estimation and identification tools.

Each of the recent approaches holds value for limited classes of problems. However, they have severe limitations when applied to the broad array of problems economists routinely confront. The danger lies in the sole reliance on these tools, which eliminates serious consideration of important policy and interpretation questions. We highlight the flexibility and adaptability of the econometric approach to causality, contrasting it with the limitations of other causal frameworks.

For instance, the NR approach does not readily accommodate unobservables and restrictions on empirical relationships produced by economic theory, which are important components of the econometric toolkit. Social interactions, peer effects, and general equilibrium theory fall outside its purview, and are currently considered frontier-topics in those fields, despite the existence of well-designed econometric tools that address these issues. These are all standard problems addressed in structural econometrics.

Similarly, the DoC approach cannot deal with the functional restrictions and covariance information routinely used in econometrics. It cannot accommodate assumptions such as monotonicity and the separability restrictions, which are essential components of modern instrumental variable analysis. The prototypical Generalized Roy model cannot be identified with DoC, although it and more general models can be identified using standard econometric tools.

This paper builds on our previous work in several ways. [Heckman \(2008a\)](#) and [Heckman and Pinto \(2015\)](#) discuss econometric causality but are less explicit than this paper

in establishing links between formal econometric models and competing approaches. We clarify the distinctions between the “do” operator of [Pearl \(2009b\)](#) and the “fix” operator of [Haavelmo \(1943\)](#) and exposit much more clearly why causality is such a difficult concept for statisticians.² We introduce a new hypothetical model that uses probabilistic tools to analyze causal models without the artifices required in competing approaches. For example, we disentangle the “SUTVA” assumption of the Neyman-Rubin model into an autonomy (structural invariance) assumption and an absence-of-general-equilibrium-effects assumption.³ We provide concrete examples of the limits and benefits of alternative causal frameworks.

This paper is organized as follows. Section 2 defines causality and discusses the tasks of causal inference. Section 3 presents the econometric model. Section 4 shows its versatility and describes various identification approaches in the Generalized Roy model. Section 5 examines the Neyman-Rubin causal model and contrasts it with the econometric approach. Section 6 investigates the *Do-Calculus* of [Pearl \(2009b\)](#). Section 7 examines non-recursive models that are ruled out in the NR approach. Section 8 summarizes the paper.

2 Causality as a Thought Experiment

A formal definition of causality relies on a modification of the same thought process used to define relationships mapping inputs X , that may contain unobserved terms, to outcomes Y using a stable map g :

$$g : X \rightarrow Y \tag{1}$$

²We note that [Haavelmo \(1943, 1944\)](#) never uses the term “fix” in his analyses. However, he introduces the notion of thought experiments in his 1943 and 1944 papers referring to them as idealized experiments. In his 1943 paper, he discusses “hypothetical splitting” of the real economic world into separate spheres of action. His example is a Keynesian consumption and investment function where he separates—hypothetically—consumer and producer actions in the context and investment are jointly determined. Haavelmo was a student of Ragnar Frisch, who defined the term econometrics and laid the foundations of econometric causal policy analysis in two foundational studies ([Frisch, 1930, 2009](#)). [Haavelmo \(1944\)](#) refers extensively to Frisch’s work and later essays on policy evaluation. He is credited with formalizing Frisch’s distinction between hypothetical worlds (models) and empirical data. See also [Bjerkholt and Dupont \(2010\)](#) and [Frisch \(2009\)](#).

³The do-calculus explicitly uses autonomous structural relationships ([Pearl, 2009b](#)).

A map is **stable** if changing its arguments over the domain of X preserves the map. Another way to express this is $Y = g(X)$, where g may be a multi-valued correspondence. An elementary version of (1) is the linear model:

$$Y = \alpha + \beta X. \tag{2}$$

In this example, stability means that α and β don't change when X is changed. This invariance property is termed *autonomy* of relationships by Frisch (1938). It is a cornerstone of causal analysis.⁴ Typical examples of autonomous relationships in economics are production functions or demand equations.

A second fundamental concept in causality is directionality. The map g states that X causes Y . Inverting this map (when possible) may produce a stable relationship, but it is, in general, not causal.

The range of Y is a set of *potential outcomes* associated with X over its domain. The map g may be either a function or a correspondence. For example, our analysis is applicable to settings such as Nash games with multiple equilibria.⁵ *Counterfactual outcomes* $Y(x)$ refer to the potential values that Y takes across different values of X . The key idea in causality is the notion captured in Alfred Marshall's phrase, "*ceteris paribus*" –all other else is equal.⁶ Comparisons of Y for different values of X – all other factors the same – are defined as *causal effects*. They are conceptual thought experiments. This definition is used explicitly in the econometric approach regardless of what is observed, the statistical properties of X and Y , the specification of functional forms for g , or how X is manipulated in any thought experiment. The Roy model (1951) is an early example of a model of two potential outcomes associated with the income that the same person would earn in different jobs. We use a generalization of it as an example prototypical model throughout this paper.

⁴Frisch (1938) defines autonomy of a function to mean functions that are "invariant" to changes in their arguments. Hurwicz (1962) prefers the term "structural" to denote autonomous equations.

⁵See e.g., Mas-Colell et al. (1995); Tamer (2003)

⁶Marshall (1961)

Issues of identification and estimation are important for making the concept of causality empirically operational, but not for defining it. However, these auxiliary issues are sometimes assumed to be paramount in defining causality in the recent non-economic literatures. For example, in an influential exposition of the Neyman-Rubin model, [Holland \(1986\)](#) insists that causal effects are only defined for experimental manipulations of X . However, issues of definition and estimation are fruitfully distinguished and are the hallmark of the econometric approach. To make our discussion more concrete, an example from the standard toolkit of empirical economics is helpful.

2.1 Regression: Conditional Expectation or Thought Experiment?

Consider the standard workhorse of empirical economics.⁷ Anticipating empirical applications, we add the distinction between observed and unobserved variables that is strictly not required for the definition of causal parameters. Consider the regression of Y on T where (Y, T) are observed and U denotes a variable that is not observed by the analyst:

$$Y = T\beta + U. \tag{3}$$

In terms of (1), $X = (T, U)$. If X is a vector of all possible causes of Y , (1) is an *all causes* model and accommodates stochastic shocks. Coupled with stability, such a model is convenient for transporting (1) to environments where different levels of T are at play (forecasting) or in combining and summarizing evidence from different studies where T varies (research synthesis).

A major source of confusion about causal models is that (3) is often defined by statisticians as a model for describing *statistical* relationships between Y and T (see e.g., [Holland, 1997](#); [Pratt and Schlaifer, 1984](#)). Doing so uses standard statistical tools to define empirical

⁷See [Haavelmo \(1943\)](#) for an early discussion of the distinction made in this section.

relationships. Note that if conditional expectations exist, $E(Y | T = t) = t\beta + E(U | T = t)$. In this approach, the statistical model could be equivalently defined as $U = Y - T\beta$.

The empirical association between T and Y operates through two channels: β and $E(U | T = t)$, unless T is mean independent of U . This approach introduces considerations about the properties of random variables that are unnecessary for defining causality in contrast to just defining an empirical regularity.

2.2 Thought Experiments

Another way to interpret $Y = T\beta + U$ is to hypothetically vary T and U : $(T, U) \rightarrow Y$ via $Y = T\beta + U$. This is not a statistical operation and lies outside standard statistics.⁸ Economists (and other scientists) use hypothetical models (thought experiments) to analyze phenomena and explore possible relationships. These and other possible relationships are not *defined* by causal operations, although they are *estimated* using statistical methods.

To clarify these ideas, it is helpful to introduce random variables ϵ_V , ϵ_T , ϵ_U which are unobserved (by the analyst) and mutually statistically independent. They are external to the model (exogeneous) and are not caused by T , U or Y .

Example 2.1. Consider four different possible causal models – all thought experiments:

Causal Model 1	Causal Model 2	Causal Model 3	Causal Model 4
$T = f_T(\epsilon_T)$	$T = f_T(\epsilon_T, \epsilon_V)$	$T = f_T(\epsilon_T, U)$	$T = f_T(\epsilon_T)$
$U = f_U(\epsilon_U)$	$U = f_U(\epsilon_U, \epsilon_V)$	$U = f_U(\epsilon_U)$	$U = f_U(\epsilon_U, T)$
$Y = T\beta + U$	$Y = T\beta + U$	$Y = T\beta + U$	$Y = T\beta + U$

In the first causal model, T does not cause U , nor does U cause T . Parameter β is the causal effect of varying T on Y for a fixed value of U . Variables T and U are statistically independent and the parameter β can be consistently estimated by OLS. In the second causal

⁸For an example of how confusing this concept is to statisticians, see [Pratt and Schlaifer \(1984\)](#) and [Holland \(1997\)](#). Holland’s confusion is significant given that he was the person who formalized the “Rubin model” ([1986](#)).

model, T does not cause U , nor does U cause T . Parameter β is still the causal effect of T on Y . However, T and U are not statistically independent because they share a common confounding variable ϵ_V and the OLS estimator of β is biased. This model is sometimes called a “common cause” model with ϵ_V being the common cause of T and U . The third causal model differs from the second model because U causes T . Although the causal relations of the second and third models differ, the causal effect of T on Y remains β . In these models, T and U are not statistically independent and the OLS estimator is generally biased.⁹ The fourth model describes the case where T causes U . In this case, the OLS estimator of the parameter β does not, in general, describes the causal effect of T on Y since we need to account for the effect of T on Y that operates through U . The OLS estimator is biased and it evaluates a combination of the direct effect of T on Y and the indirect effect of T on Y via U .

Using the standard regression model as a starting point blurs the logic of this thought process. Econometrics textbooks commonly introduce causality in the context of the linear model (3). In this approach, the identification of causal effects is often reduced to a statistical property of the econometric model, namely, that causal effects can be assessed when variables T and U are uncorrelated. It gives rise to the practice of defining causal effects as conditional probability statements instead of statements about manipulating variables in a thought experiment.

In fact, OLS is based on statistical assumptions that are void of any causal interpretation. The OLS fitted value for the outcome Y conditioning on $T = t$ evaluates the conditional expectation $E(Y | T = t)$ instead of the counterfactual expectation $E(Y(t) | T = t)$, where the counterfactual outcome $Y(t)$ is the value of Y when T is externally set to a value t . The causal content of the OLS model arises only when we invoke concepts such as fixing and counterfactuals. These concepts are not part of the standard statistical toolkit. Whether or

⁹Thus, $Y(t) \perp\!\!\!\perp T|U$ holds for the third model but not for the second model.

not we can identify β in a sample is an entirely separate question from defining the causal impact of T on Y .

Frisch, the founding father of modern econometric causal policy analysis, clearly understood that the study of causality is an exercise in abstract thought, and that “*Causality is in the Mind*”:

“... we think of a cause as something imperative which exists in the **exterior world**. In my opinion this is fundamentally **wrong**. If we strip the word cause of its animistic mystery, and leave only the part that science can accept, nothing is left except a certain way of thinking. [T]he scientific ... problem of **causality** is essentially a problem regarding our **way of thinking**, not a problem regarding the nature of the exterior world.” — [Frisch \(1930\)](#), p. 36

Stated differently, Frisch is saying causality is the outcome of a thought experiment, i.e., a model.

2.3 The Econometric Approach to Causality

The econometric approach to causality develops explicit hypothetical models where inputs cause outcomes. A common context is the study of policy evaluations when economic agents choose treatments that affect economic outcomes of interest. “Treatments” are inputs (the T) which need not be restricted to binary or discrete valued variables. The mechanisms governing the choice of inputs is central to the study of the causal effect of treatment on outcomes. Identification/estimation/interpretation of empirical counterparts to the hypothetical counterfactuals require careful accounting of unobserved (by the analyst) variables (U) that cause both input choice and outcomes. Structural econometric models do just that.¹⁰

¹⁰Caricatures sometimes made in the non-economic literatures that the choices of inputs T involve highly stylized rational choice models or perfect information are false (see, e.g., [Morgan and Winship, 2015](#)). Some hypothetical models might maintain those assumptions, but such assumptions are in no way essential to the enterprise.

2.4 Four Distinct Policy Questions

The econometric approach to causality distinguishes four distinct classes of policy problems and addresses each of them, sometimes in the same analysis.¹¹

P1 *Evaluating the impacts of implemented interventions on outcomes in a given environment, including their impacts in terms of the well-being of the treated and society at large. The simplest forms of this problem are typically addressed in the non-economic literatures: does a program in place “work” in terms of policy impacts?*

The non-economic literatures addressing **P1** identify and estimate treatment effects (most often average treatment effects) without investigating how they arise or whether alternative programs might be better or even what “better” means. In terms of our linear equation example, it seeks to know the sign and magnitude of β . However, most economic and policy analysts seek greater generality for their findings. This leads to problem **P2**.

P2 *Understanding the mechanisms producing treatment effects and policy outcomes.*

This asks the analyst to investigate the causes of effects and is a central task of economic theory and policy analysis.¹² It embeds (3) in a model that explains how T operates (i.e., which factors explain the $Y - T$ relationship). It goes beyond the coarse description of “treatment” T to explicate the factors that produce Y . It links with **P3** and **P4** stated below to consider how alternative mechanisms generate observed outcomes and can be used to forecast policies going forward, or explain the findings of any given study in a particular environment. **P2** is also an integral part of the task of constructing alternatives to maintained hypotheses and interpreting evidence using economic models.

¹¹See Heckman (2008a).

¹²Holland (1986) features the narrow goal of investigating the “effects of causes” in his definition of the Neyman-Rubin model.

P3 *Forecasting the impacts (constructing counterfactual states) of interventions implemented under one environment when the intervention is applied to other environments, including their impacts in terms of well-being.*

This goes beyond **P2** to interpret why outcomes vary among environments. It replaces crude meta-analysis of treatment effects with principled explanations of mechanisms and their impacts and extrapolates mechanisms to other environments to answer **P1** in those environments.¹³ Structural models are useful vehicles for summarizing evidence from multiple studies.¹⁴ Forecasting in new environments is a traditional problem in econometrics (see, e.g., [Theil, 1958](#); [Hamilton, 2000](#); [Chatfield, 2000](#)). However, the truly ambitious problem addressed by policy analysts is **P4**.

P4 *Forecasting the impacts of interventions (constructing counterfactual states associated with interventions) never previously implemented to various environments, including their impacts in terms of well-being.*

This is a fundamental challenge addressed in econometric policy analysis. This problem motivated the creation of econometric causal models.¹⁵ It is also a central feature of the scientific analysis of empirical regularities.

One impetus for the econometric structural approach was to conduct policy analysis for the post-World War II era using models fit on data from the pre-World War II Depression era. Econometric policy analysis is the vehicle for framing and addressing the likely impacts of new policies and new environments, never previously experienced. [Marschak \(1953\)](#) provides an insightful discussion of this task in the context of forecasting the impact of new economic policies using data collected in environments in which the proposed policies were not in place.¹⁶ The often-cited “critique” of [Lucas \(1976\)](#) updates Marschak’s policy analysis to

¹³Recent work in computer science has begun to reinvent the logic of econometric forecasting using its own colorful private language but without any fresh insights or acknowledgement of a large body of econometric thought (see, e.g., [Bareinboim and Pearl, 2016](#)).

¹⁴See, e.g., [Burszтын and Yang \(2021\)](#) or [Nerlove \(1967\)](#).

¹⁵See [Frisch \(1930, 1933, 1938\)](#) and [Tinbergen \(1930\)](#).

¹⁶[Knight \(1921\)](#) succinctly states the problem and its solution in his enigmatic remark, “*the existence of a problem of knowledge depends on the future being different from the past, while the possibility of a solution of*

stochastic environments. [McFadden \(1974\)](#) is a Nobel-Prize winning example of how a leading economist who successfully met this challenge in forecasting the demand for a new transportation system in the San Francisco Bay area.

The econometric approach distinguishes three tasks of econometric causal policy analysis that are often conflated in the non-economic statistical literatures:

Table 1: **Three Distinct Tasks in Econometric Causal Analysis**

Task	Description	Requirements	Types of Analysis
1: Model Creation	Defining the class of hypotheticals or counterfactuals by thought experiments (models)	A scientific theory: A purely mental activity	Outside Statistics; Hypothetical Worlds
2: Identification	Identifying causal parameters from hypothetical populations	Mathematical analysis of point or set identification; this is a purely mental activity	Probability Theory
3: Estimation	Estimating parameters from real data	Estimation and testing theory	Statistical Analysis

Our regression example illustrates these distinctions. Models for counterfactuals do not necessarily require any statistical analysis. Identification is a separate issue required to recover β from hypothetical model distributions of data where statistical variation is not an issue.¹⁷ Estimation, on the other hand, considers how to recover model parameters from empirical sampling distributions where statistical variation is a concern. Trygve Haavelmo, a student of Frisch, developed an empirically operational econometric framework for causal analysis that distinguished these three tasks ([1943](#); [1944](#)).

the problem depends on the future being like the past.” Knight meant that analysts use ingredients estimated on historical data to construct forecasts of the unknown. This is a task that involves judgements and insights about invariant mechanisms beyond straight applications of fitted statistical models.

¹⁷[Lewbel \(2019\)](#) and [Fisher \(1966a\)](#) are definitive treatments of identification in economics.

3 Econometric Causal Models

Econometric causal models are flexible frameworks that can be used to address a variety of economic policy problems that are not naturally squeezed into simple “treatment-control” frameworks. They go well beyond the narrow treatment effect literature to address the following topics listed in Table 2¹⁸ :

Table 2: Problems Addressed by Econometrics

- (a) Investigate the causes of effects, not just the effects of causes – the goal of the treatment effect literature announced by Holland (1986) in defining the “Rubin model;”
- (b) Interpret empirical relationships within economic choice and outcome frameworks;
- (c) Analyze data using *a priori* information from theory;
- (d) Account systematically for shocks, errors by agents, and measurement errors;
- (e) Analyze dynamic models;
- (f) Accommodate multiple approaches to identification beyond randomization, instrumental variables, and matching that exploit restrictions within and across equations on causal relationships produced by theory;
- (g) Exploit restrictions across equations and unobservables within and across equations to identify causal parameters;
- (h) Make forecasts in new environments;
- (i) Synthesize evidence across studies using common parameters embedded in common conceptual frameworks rather than crude statistical meta-analysis;
- (j) Make forecasts of new policies never previously implemented; and
- (k) Analyze interactions across agents within markets and also within social settings (general equilibrium and peer effects).

Econometric methodology for establishing causality is comprehensive and adaptable, as it is specifically designed to address a wide array of causal questions pertinent to economics. In

¹⁸Table 2 is only a partial list of the rich array of problems addressed by the econometric approach.

contrast, alternative causal frameworks are often not conceived with the specific investigative needs of economists in mind. Consequently, these methods are typically tailored to address only a specific subset of causal questions, primarily focusing on the application of a limited range of techniques to specialized categories of problems, predominantly those within the problem class P1.

Alternative methodologies, such as the NR approach, can be highly effective in analyzing causal effects such as average treatment effects or the effect of treatment on the treated within the contexts of Randomized Controlled Trials (RCTs). However, their utility becomes markedly constrained when addressing the more complex causal questions mentioned in Table 2. This is a consequence of *Marschak's Maxim* (Heckman, 2008a) that for certain narrowly focused problems, special versions of the econometric approach are highly effective. One need not necessarily implement more general models that address a wider set of questions when addressing specific focused problems. However, such models are by design, of limited value in addressing wider classes of problems. We now state the econometric model formally using the convenient tool of graph theory that is widely used in many branches of applied mathematics.

3.1 Econometric Causal Models

Heckman and Pinto (2015) develop a causal framework that formalizes Frisch's insight that causality is the outcome of a thought experiment and places Haavelmo's approach (1943; 1944) in the framework of more recent policy evaluation models. They distinguish an *empirical model* that generates the observed data from a *hypothetical model* that formalizes the thought experiments of manipulating inputs that defining causality. The empirical model describes the data generating process, which differs from the hypothetical model which is an abstract model that is a thought experiment. They place the definition and operationalization of causality in a probabilistically consistent approach that does not require special

rules or procedures invented to characterize causality that are essential features of some of the non-economic literature.

Some notation is useful in describing the framework. We borrow it from the literature in applied mathematics. Dawid (1979) is a major source of conditional independence relationships. Lauritzen (1996) is a concise treatment of the graph theory we use.

3.2 A Causal Model

A causal model \mathbb{M} is a system of policy-invariant (autonomous) structural equations like (1) that characterize the mapping $\mathbb{M} : \mathcal{T} \rightarrow \mathbb{P}(\mathcal{T})$ between a set of variables \mathcal{T} and its power set $\mathbb{P}(\mathcal{T})$. Elements in \mathcal{T} are random variables or random vectors that may be observed or unobserved by the analyst. It is convenient to define the set $\mathcal{E} = \{\epsilon_K; K \in \mathcal{T}\}$ which contains an error term ϵ_K for each $K \in \mathcal{T}$. Error term ϵ_K shares the same dimension as K . This variable is assumed present even if there are additional unobserved variables. This formal device allows us to avoid degenerate random variables so standard tools of probability theory can be used.

A structural equation for a variable $K \in \mathcal{T}$ is an autonomous function denoted by $f_K : (\mathbb{M}(K), \epsilon_K) \rightarrow \mathbb{R}^{|K|}$. Variables in $\mathbb{M}(K)$ are said to directly cause K . In recursive formulations, a variable cannot directly cause itself, that is, $K \notin \mathbb{M}(K)$ for all $K \in \mathcal{T}$. We relax recursivity in a later section, where we discuss simultaneous equation models in which sets of outcome variables are jointly determined.

A variable T not caused by any variable, so $\mathbb{M}(T) = \emptyset$, is called *external*. In this case, its structural function is given by $T = f_T(\epsilon_T)$. Error terms are externally-specified (i.e., exogenous). This means that error terms in set \mathcal{E} are not caused by any variable in \mathcal{T} . We impose, without loss of generality, that error terms are mutually statistically independent.¹⁹

¹⁹The independence among error terms comes without loss of generality as any dependence structure could be modeled via other unobserved variables in \mathcal{T} .

All variables are defined on a common probability space $(\mathcal{I}, \mathcal{F}, P)$, using standard notation for σ - algebras.

3.3 The Generalized Roy Model

We use the Generalized Roy model as our leading example of a structural model. It is a cornerstone of the literature in applied economics and policy evaluation.²⁰ The original Roy model of counterfactuals (1951) analyzed earnings inequality in two sectors of the economy. All persons have two potential incomes: $Y(0)$ in Sector 0 and $Y(1)$ in Sector 1. Agents choose sectors based on their perceived net benefit I . In the simplest case, the benefit is the income gain $I = Y(1) - Y(0)$. More general models allow for costs, like tuition, migration costs, and psychic costs of participation. Potential incomes $(Y(0), Y(1))$ depend on observed variables X while benefit I may depend on $Y(1), Y(0)$, and X and an externally specified variable Z , which may be policy variables that influences participation costs. The agent's choice of sector is given by $T = \mathbf{1}[I(X, Z) > 0]$. The model has been generalized to analyze multiple sectors and dynamic discrete choices (see Abbring and Heckman, 2007; Heckman and Vytlacil, 2007a,b).

The most common form of the model is:

$$Y(1) = g_1(X) + U_1 \tag{4}$$

$$Y(0) = g_0(X) + U_0 \tag{5}$$

$$C = g_c(Z, X) + U_c \tag{6}$$

$$T = \mathbf{1}[Y(1) - Y(0) - C \geq 0]. \tag{7}$$

g_1, g_0 , and g_c are autonomous functions. The variables X are observed and cause the outcome and choices. Variable Z serves as an instrumental variable. It is not an argument of the outcome equations. Variables U_1, U_0 , and U_c are exogenous and unobserved variables

²⁰See, e.g., Heckman and Taber (2008); Heckman and Vytlacil (2007a,b).

that are statistically independent of Z, X , namely, $(U_1, U_0, U_c) \perp\!\!\!\perp (X, Z)$.²¹ Choice theory as embodied in (7) helps in determining relevant variables Z which can serve as instrumental variables.

The individual level treatment effect is $Y(1) - Y(0)$. The evaluation problem arises because for each person we observe either $Y(0)$ or $Y(1)$, but not both. We observe $Y(1)$ if $T = 1$ and $Y(0)$ if $T = 0$, namely $Y = T \cdot Y(1) + (1 - T) \cdot Y(0)$.²²

Z affects Y only through its influence on T . The typical analysis reformulates the analysis at the population level rather than at the individual level. A common parameter of interest is the average treatment effect $ATE = E(Y(1) - Y(0))$ which is the mean treatment effect across all agents. Treatment on the treated focuses on $TOT = E(Y(1) - Y(0)|T = 1)$. The probability distribution of the counterfactual outcomes $Y(t); t \in \{0, 1\}$ are sometimes investigated.

The Generalized Roy model has been extended in many ways.²³ The model is systematically ignored in the non-economic literatures, despite its intellectual priority and relevance.²⁴ The Generalized Roy model allows for multiple choices. It can account for subjective evaluations of the benefits of each choice by subsuming variables U_1, U_0, U_c , in an unobserved random vector V that causes both T and Y (see Heckman and Pinto, 2018; Heckman and Vytlacil, 2007a,b).

A simple yet general representation of the Generalized Roy model comprises four random variables $\mathcal{T} = \{Z, V, T, Y\}$, where Z is an instrumental variable that causes an outcome Y only through its effects on a treatment choice T . The variable V denotes an externally-

²¹This independence relationship may also take the form of the conditional independence $(U_1, U_0, U_c) \perp\!\!\!\perp Z|X$.

²²This switching regression relationship was first used by Quandt (1958). See also Quandt (1988).

²³For instance, Heckman and Vytlacil (2007a) investigate multiple versions of the original model. Heckman et al. (2008) extend the model to ordered and general unordered choice models. Heckman and Pinto (2018) and Lee and Salanié (2018) investigate the case of unordered multiple choice models with multi-valued treatments. Abbring and Heckman (2007) consider dynamic discrete choice models in this framework.

²⁴See e.g., Holland (1986); Imbens and Rubin (2015); Pearl (2009b, 2012); Rubin (1974, 1978).

specified (exogenous) and unobserved confounding variable that causes both T and Y .²⁵ In the context of the Generalized Roy model, Z stands for external policy or “shifter” vectors. V is a source of selection bias as it induces covariation between choice T and outcome Y that is not due to the causal effect of the treatment T on the outcome Y . For now, we suppress the X variables for the sake of notational simplicity. Table 3 displays four equivalent representations of the Generalized Roy model.

Table 3: **Representations of the Generalized Roy Model**

	Variable Map	Structural Eq.	DAG	LMC
Z	$\mathbf{M}(Z) = \emptyset$	$Z = f_Z(\epsilon_Z)$	<pre> graph LR Z[Z] --> T[T] T --> Y[Y] V((V)) --> T V --> Y </pre>	$Z \perp\!\!\!\perp V \emptyset$
V	$\mathbf{M}(V) = \emptyset$	$V = f_V(\epsilon_V)$		$V \perp\!\!\!\perp Z \emptyset$
T	$\mathbf{M}(T) = \{Z, V\}$	$T = f_T(Z, V, \epsilon_T)$		$T \perp\!\!\!\perp \emptyset (Z, V)$
Y	$\mathbf{M}(Y) = \{T, V\}$	$Y = f_Y(T, V, \epsilon_Y)$		$Y \perp\!\!\!\perp Z (T, V)$

The first column of Table 3 lists the variables of the Roy model. The second column describes the causal model as a mapping of the variable set. The third column displays the corresponding structural equations. The fourth column displays the model as a Directed Acyclic Graph (DAG), where arrows denote causal relationships, circles denote unobserved variables, and squares denote observed variables.²⁶ To avoid clutter, we keep the ϵ implicit.

The last representation of Table 3 uses a property called the Local Markov Condition (LMC).²⁷ Some notation is necessary to state the condition. The language of Bayesian networks uses the term parents of K for the variables that directly cause K , that is $\mathbf{M}(K)$. Children of K comprise the variables directly caused by K , namely, $\text{Ch}(K) = \{J \in \mathcal{T}; K \in \mathbf{M}(J)\}$. The descendants of a variable K , $\text{D}(K)$, include all variables that are directly or

²⁵Choice T may be binary, discrete or continuous and the confounder variable V can denote a random vector of arbitrary dimension.

²⁶We refer to Lauritzen (1996) for background on DAGs and Bayesian Networks.

²⁷See Kiiveri et al. (1984); Pearl (1988) for further information on the Local Markov Condition.

indirectly caused by K . These include all the subsequent iterations of the children of K .²⁸ A causal model is recursive (acyclic) if no variable is a descendant of itself.

The LMC is a property of recursive models stating that a variable is independent of its non-descendants conditioned on its parents.

$$\mathbf{LMC:} \quad K \perp\!\!\!\perp (\mathcal{T} \setminus \mathbb{D}(K)) \mid \mathbb{M}(K) \quad (8)$$

For instance, outcome Y has no descendants and its parents are $\{V, T\}$. Thus its LMC is $Y \perp\!\!\!\perp Z \mid (T, V)$, as listed in the last row of Table 3. Z has no parents and its descendants are T, Y . The set of LMC for all variables in \mathcal{T} fully characterizes the causal model. Additional independence relationships may be generated by the Graphoid Axioms²⁹ of Dawid (1976) or through graphical methods such as the d -separation criteria of Geiger et al. (1990).

3.4 Counterfactual Approaches: Formalizing Frisch’s Insight

Frisch’s statement that “Causality is in the Mind” means that the causal analysis of treatment T relies on a thought experiment that assigns values to the treatment variable in a fashion external to the system analyzed. This hypothetical manipulation of T affects only the variables caused by T . Specifically, changing T affects its descendant Y but not its ancestors V, Z .

²⁸Notationally, for any subset $\tilde{\mathcal{T}} \subset \mathcal{T}$, let $\text{Ch}(\tilde{\mathcal{T}})$ be the union of the children of all the variables in $\tilde{\mathcal{T}}$, that is, $\text{Ch}(\tilde{\mathcal{T}}) = \cup_{K \in \tilde{\mathcal{T}}} \text{Ch}(K)$. The descendants of K is the smallest set $\mathbb{D}(K) \subset \mathcal{T}$ that contains the children of K , $\text{Ch}(K) \subset \mathbb{D}(K)$, and its own children, $\text{Ch}(\mathbb{D}(K)) = \mathbb{D}(K)$.

²⁹Dawid (1976) defines Graphoid Axioms consist of six rules that apply for any disjoint sets of variables $X, W, Z, Y \subseteq \mathcal{T}$:

(A)Symmetry:	$X \perp\!\!\!\perp Y \mid Z \Rightarrow Y \perp\!\!\!\perp X \mid Z.$
(B)Decomposition:	$X \perp\!\!\!\perp (W, Y) \mid Z \Rightarrow X \perp\!\!\!\perp Y \mid Z.$
(C)Weak Union:	$X \perp\!\!\!\perp (W, Y) \mid Z \Rightarrow X \perp\!\!\!\perp Y \mid (W, Z).$
(D)Contraction:	$X \perp\!\!\!\perp W \mid (Y, Z) \text{ and } X \perp\!\!\!\perp Y \mid Z \Rightarrow X \perp\!\!\!\perp (W, Y) \mid Z.$
(E)Intersection:	$X \perp\!\!\!\perp W \mid (Y, Z) \text{ and } X \perp\!\!\!\perp Y \mid (W, Z) \Rightarrow X \perp\!\!\!\perp (W, Y) \mid Z.$
(F)Redundancy:	$X \perp\!\!\!\perp Y \mid Z \Rightarrow X \perp\!\!\!\perp Y \mid Z.$

Frisch’s thought experiment is conceptually simple. However, it is a causal operation outside the scope of statistical theory. In statistics, random variables are fully characterized by their joint distributions. This information by itself is insufficient for causal analysis as it lacks directionality – a central feature of causal models. Frisch’s thought experiment uses additional information on causal direction when it partitions the variables studied into those caused by T and those that are not.

Frisch’s thought experiment was formalized through the use of the “*fix*” or “set” operator implicit in the seminal work of [Haavelmo \(1943\)](#). *Counterfactual* outcomes are obtained by the hypothetical (external) manipulation of the targeted variable that causes the outcome of interest. In the Roy model, the counterfactual outcome $Y(t)$ is obtained by *fixing* the T -input of the outcome equation to a value $t \in \text{supp}(T)$ so that $Y(t) = f_Y(t, V, \epsilon_Y)$. Fixing only affects the outcome equation. It substitutes the treatment random variable T by the treatment value t . It makes all descendants of T functions of the fixed value of $T = t$. It does not eliminate the equation for T from the causal model nor does it modify the choice equation $T = f_T(Z, V, \epsilon_T)$.

The *do-operator* of [Pearl \(1995, 2012\)](#) operates in a fashion similar to fixing as it substitutes all T -inputs from structural equations of the variables directly caused by T . The *do-operator* differs from fixing by *deleting* (“shutting down”) the structural equation for the treatment variable T , which effectively suppresses the determining equation of random variable T from the causal model and replaces it with a fixed value $T = t$ that affects all descendent relationships. Eliminating this equation excludes the possibility of defining parameters like TOT that condition on T .

Neither *fix* nor *do* are well-defined in statistics. They are causal operators that only affect the distribution of the descendants of the variable being fixed. In contrast, statistical conditioning affects the distributions of all variables that are not statistically independent of the conditioning variable. Fixing T in the Roy model affects the outcome Y but does not impact the confounder V or the instrument Z , which remain statistically independent.

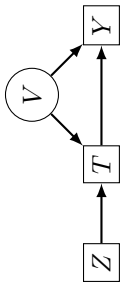
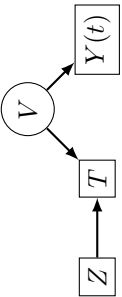
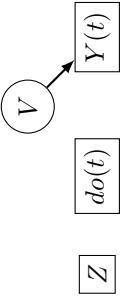
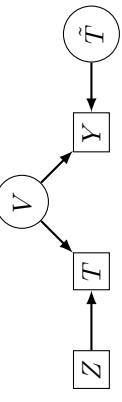
Conditioning on T , on the other hand, alters the distributions of Z and V , which are no longer statistically independent.

Heckman and Pinto (2015) develop a causal framework that expresses the causal operations of *fixing* or *doing* in a framework using standard statistical tools. They distinguish the *empirical model* that generates observable data from a *hypothetical model* that is used to formulate the thought experiments involving the manipulation of inputs determining causality. The hypothetical model is an abstract model (thought experiment) that shares the same structural equations and the same distributions of error terms as the empirical model. It differs from the empirical model by appending a hypothetical variable \tilde{T} that replaces the T -input affecting descendants of T . The hypothetical variable captures the causal notion of an external manipulation of treatment. The hypothetical model operates downstream of T and translates the causal operation of fixing T into the statistical operation of conditioning on \tilde{T} .

Notationally, we use $\mathcal{T}_e, \mathcal{E}_e, \mathbb{M}_e, P_e, E_e$ for the variable set, error terms, causal model, probability, and expectation of the empirical model, while $\mathcal{T}_h, \mathcal{E}_h, \mathbb{M}_h, P_h, E_h$ denote the counterparts in the hypothetical model. The hypothetical model replaces all T -inputs. In this case, the hypothetical and empirical model are related in the following fashion: (1) the hypothetical model has an additional variable \tilde{T} , $\mathcal{T}_h = \mathcal{T}_e \cup \{\tilde{T}\}$; (2) the hypothetical variable causes all descendants of T , $\mathbb{M}_h(K) = (\mathbb{M}_e(K) \cup \{\tilde{T}\}) \setminus \{T\}$ for all $K \in \mathbb{D}_e(T)$; (3) variable T has no descendants in the hypothetical model, that is, $\mathbb{D}_h(T) = \emptyset$; (4) all remaining causal relations stay the same, that is, $\mathbb{M}_h(K) = \mathbb{M}_e(K)$ for all $K \in \mathcal{T}_e \setminus \{\mathbb{D}_e(T) \cup \{T\}\}$.

It is useful to illustrate these ideas using the Generalized Roy Model. For notational clarity, we use \mathbb{M}_e for the empirical (original) model, \mathbb{M}_{fix} for the model that applies the *fix*-operator, \mathbb{M}_{do} for the *do*-operator, and \mathbb{M}_h for the hypothetical model. We also use the subscripts e, fix, do, h for the probability distributions, expectations associated to each model. Table 4 displays the Roy model for each of these frameworks.

Table 4: Generalized Roy Model: Approaches to Generating Counterfactuals

Empirical Models		Hypothetical Model	
Empirical Model (\mathbb{M}_e)	Fixing T at t (\mathbb{M}_{fix})	Doing $do(t)$ (\mathbb{M}_{do})	Hypothetical Model (\mathbb{M}_h)
<i>Structural Equations</i>			
$ \begin{aligned} V &= f_V(\epsilon_V) \\ Z &= f_Z(\epsilon_Z) \\ T &= f_T(Z, V, \epsilon_T) \\ Y &= f_Y(T, V, \epsilon_Y) \end{aligned} $ 	$ \begin{aligned} V &= f_V(\epsilon_V) \\ Z &= f_Z(\epsilon_Z) \\ T &= f_T(Z, V, \epsilon_T) \\ Y(t) &= f_Y(t, V, \epsilon_Y) \end{aligned} $ 	$ \begin{aligned} V &= f_V(\epsilon_V) \\ Z &= f_Z(\epsilon_Z) \\ do(T=t) & \\ Y(t) &= f_Y(t, V, \epsilon_Y) \end{aligned} $ 	$ \begin{aligned} V &= f_V(\epsilon_V) \\ Z &= f_Z(\epsilon_Z) \\ T &= f_T(Z, V, \epsilon_T) \\ Y &= f_Y(\tilde{T}, V, \epsilon_Y) \\ \tilde{T} &= f_{\tilde{T}}(\epsilon_{\tilde{T}}) \end{aligned} $ 
<i>Directed Acyclic Graphs (DAGs)</i>			
<i>Local Markov Conditions</i>			
$ \begin{aligned} V &\perp\!\!\!\perp Z \\ Z &\perp\!\!\!\perp V \\ T &\perp\!\!\!\perp \emptyset \mid (Z, V) \\ Y &\perp\!\!\!\perp Z \mid (T, V) \end{aligned} $ <p>(not defined for the model)</p>	$ \begin{aligned} V &\perp\!\!\!\perp Z \\ Z &\perp\!\!\!\perp (V, Y(t)) \\ T &\perp\!\!\!\perp Y(t) \mid (Z, V) \\ Y(t) &\perp\!\!\!\perp (Z, T) \mid V \end{aligned} $ <p>(not defined for the model)</p>	$ \begin{aligned} V &\perp\!\!\!\perp Z \\ Z &\perp\!\!\!\perp (V, Y(t)) \\ Y(t) &\perp\!\!\!\perp Z \mid V \end{aligned} $ <p>(not defined for the model)</p>	$ \begin{aligned} V &\perp\!\!\!\perp (Z, \tilde{T}) \\ Z &\perp\!\!\!\perp (V, Y, \tilde{T}) \\ T &\perp\!\!\!\perp (\tilde{T}, Y) \mid (Z, V) \\ Y &\perp\!\!\!\perp (Z, T) \mid (\tilde{T}, V) \\ \tilde{T} &\perp\!\!\!\perp (T, V, Z) \end{aligned} $
<i>Factorial Decomposition of the Joint Probability Distributions</i>			
$ P_e(Y T, V)P_e(T Z, V)P_e(V)P_e(Z) $	$ P_{fix}(Y(t), T, V, Z) = P_{fix}(Y(t) V)P_{fix}(T V, Z)P_{fix}(V)P_{fix}(Z) $	$ P_{do}(Y(t), V, Z) = P_{do}(Y(t) V)P_{do}(V)P_{do}(Z) $	$ P_h(Z, V, T, \tilde{T}, Y) = P_h(Y \tilde{T}, V)P_h(T Z, V)P_h(V)P_h(Z)P_h(\tilde{T}) $

Subscript e denotes empirical (original) model. Subscript fix denotes the model that uses the fix operator, that is when treatment T is fixed to t . Subscript do denotes the model that employs the do-operator. Subscript h denotes the hypothetical model. Notice further that Z and V are externally specified so that $P_e(Z) = P_{fix}(Z) = P_{do}(Z) = P_h(Z)$ and $P_e(V) = P_{fix}(V) = P_{do}(V) = P_h(V)$.

The first column of Table 4 presents the original empirical model. The second and third columns present the models generated by the *fix* and the *do* operators respectively. Both models constraint the T -input of the outcome equation by a value $t \in \text{supp}(T)$. The main difference between these models is that *fix* retains the equation for treatment while *do* suppresses it. The hypothetical model is displayed in the last column of Table 4. It replaces the T -input of the outcome equation with an external hypothetical variable \tilde{T} .

The first panel presents the structural equations of each approach. The second panel display the models as DAGs. The third panel describes the independence relationships generated by each causal model, and the last panel of the table presents the factorization of the joint distribution of the model variables. We use P_e for the probability distribution of the empirical model, P_{fix} for the model generated by the fix operator, P_{do} for the *do* operator and P_h for the hypothetical model. The factorizations differ according to the number of variables and causal relations of each counterfactual model.

The empirical (\mathbf{M}_e), *fix* (\mathbf{M}_{fix}), and hypothetical (\mathbf{M}_h) models share the same distributions of error terms $\epsilon_Z, \epsilon_V, \epsilon_T, \epsilon_Y$. Therefore the joint distribution of non-descendant T , that is (V, Z) , is the same across these models. The *do* model eliminates the error term ϵ_T , and the distribution of T is not defined.

The structural equation for the counterfactual outcome $Y(t)$ in the *fix* or *do* models depends only on V and ϵ_Y and thus the models have the same distribution of $Y(t)$. The hypothetical variable \tilde{T} enables us to circumvent the necessity of introducing a special causal operator. The variable has no parents and, according to the LMC (8), it is independent of all its non-descendants, $\tilde{T} \perp\!\!\!\perp (T, V, Z)$. In particular, $\tilde{T} \perp\!\!\!\perp T$ always hold for any hypothetical model \mathbf{M}_h . \tilde{T} is also statistically independent of error terms $\epsilon_Z, \epsilon_T, \epsilon_V, \epsilon_Y$, and the counterfactual outcome is obtained by simply conditioning on \tilde{T} . In summary, we have that:

$$\left(Y \mid \tilde{T} = t\right)_{\mathbf{M}_h} \stackrel{d}{=} \left(Y(t)\right)_{\mathbf{M}_{fix}} \stackrel{d}{=} \left(Y(t)\right)_{\mathbf{M}_{do}}. \quad (9)$$

It is also the case that equation (9) holds when conditioned on any variable K that is non-descendant variable of \tilde{T} , namely, Z, V and T .

To fix ideas, let T be an indicator of college graduation and Y denote adult income. Treatment-on-the-treated (TOT) is the average causal effect of college on income for those who choose to go to college ($T = 1$), which is $TOT = E_{fix}(Y(1) - Y(0) | T = 1)$ using the fix operator. The parameter is equivalently described as $TOT = E_h(Y | \tilde{T} = 1, T = 1) - E_h(Y | \tilde{T} = 0, T = 1)$ using the hypothetical model. The *do* operator *excludes* the treatment variable T , which poses a challenge in defining the TOT parameter. [Shpitser and Pearl \(2009\)](#) solve this issue by adding additional special structure to their counterfactual model.

Equation (9) may suggest that the way that counterfactuals are expressed is of little relevance in the study of causality. That assessment is quite misleading. Small differences in characterizing counterfactuals have significant consequences for the machinery used to identify causal effects. Section 6 illustrates the difference between an identification analysis using the *do*-calculus and an identification analysis using the hypothetical model framework. Section 5 compares identification in NR with identification in the structural model.

3.5 Identification of Counterfactual Outcomes

We next consider Task 2 in Table 1. Counterfactuals are said to be identified if they can be expressed in terms of the observed data generated by the empirical model M_e . This task requires us to connect the probability distribution (or expectation) of counterfactual variables with the population distributions of the empirical model. The mechanics for establishing this connection depends on which causal model is used to describe counterfactuals.

First consider the *fix* operator of model M_{fix} in Table 4. The LMC of $Y(t)$ in M_{fix} implies that:

$$Y(t) \perp\!\!\!\perp T | V. \tag{10}$$

Equation (10) states that the counterfactual outcome $Y(t)$ is independent of the treatment variable T conditional on the confounding variable V . This relationship is an example of a *matching condition*. It helps identify treatment effects as it connects the counterfactual outcome $Y(t)$ in \mathbb{M}_{fix} with the empirical model \mathbb{M}_e :

$$P_{fix}(Y(t) | V) = P_{fix}(Y(t) | V, T = t), \quad (11)$$

$$= P_{fix} \left(\sum_{t \in \text{supp}(T)} \mathbf{1}[T = t] Y(t) | V, T = t \right), \quad (12)$$

$$= P_{fix} \left(\sum_{t \in \text{supp}(T)} \mathbf{1}[T = t] f_Y(t, V, \epsilon_Y) | V, T = t \right), \quad (13)$$

$$= P_{fix}(f_Y(T, V, \epsilon_Y) | V, T = t), \quad (14)$$

$$= P_e(Y | V, T = t). \quad (15)$$

Equations (11)–(15) use structural equations to express the probability distribution of the counterfactual outcome $Y(t)$ in \mathbb{M}_{fix} with the distribution of the outcome Y in empirical model \mathbb{M}_e . The first equation (11) is due to the matching condition (10). Equations (11)–(14) apply the definition of the structural equations. The last equation (15) uses the fact that variables T, V, ϵ_Y share the same distribution in both models \mathbb{M}_{fix} and \mathbb{M}_e .

The hypothetical model \mathbb{M}_h offers criteria that enable analysts to connect the counterfactual and empirical distributions in a systematic manner. For any disjoint set of variables Y, W in $\mathcal{T}_h \setminus \{T, \tilde{T}\}$ and any values $t, t' \in \text{supp}(T)$ we have that:³⁰

$$Y \perp\!\!\!\perp \tilde{T} | (T, W) \Rightarrow P_h(Y | \tilde{T} = t, T = t', W) = P_h(Y | T = t', W) = P_e(Y | T = t', W), \quad (16)$$

$$Y \perp\!\!\!\perp T | (\tilde{T}, W) \Rightarrow P_h(Y | \tilde{T} = t, T = t', W) = P_h(Y | \tilde{T} = t, W) = P_e(Y | T = t, W). \quad (17)$$

Equations (16)–(17) state two conditions that involve independence relationships in the hypothetical model. They state that we can switch from the hypothetical to the empiri-

³⁰See Heckman and Pinto (2015) for a proof. The criteria (16)–(17) still holds if the values $t, t' \in \text{supp}(T)$ were replaced by subsets $\mathcal{A}, \mathcal{A}' \subset \text{supp}(T)$ respectively.

cal model whenever the hypothetical model yields the independence relationships (16) and (17).³¹

The application of these rules is simple. For example, the LMC of Y in M_h of Table 4 generates the following matching condition:

$$Y \perp\!\!\!\perp T | (\tilde{T}, V). \quad (18)$$

Thus, according to (17), we have that $P_h(Y | \tilde{T} = t, V) = P_e(Y | T = t, V)$.

The hypothetical framework gives a systemic approach for connecting hypothetical and empirical models. The framework employs additional structure beyond what is obtained from fixing that might not be required in analyzing the simple Roy model. Section 6 explores more complex models where the additional complexity of the hypothetical framework is warranted.

The *do* operator does not generate a matching condition such as (10) or (18) because the equation for treatment T is absent. Instead, the do-calculus of (Pearl, 2009b) checks for matching conditions using a DAG-based analysis called the “back-door” criterion Pearl (1993). The method employs special jargon that may be obscure to most economists. The criterion is part of the do-calculus, which consists of a set of DAG-oriented techniques that enables us to systematically examine the identification of causal effects. The method is general in the sense that it applies to any DAG, but limited in the sense that it does not accept identifying assumptions outside the DAG terminology. We discuss the do-calculus machinery, its benefits and limitations in Section 6.

The counterfactual models M_{fix} , M_h and M_{do} employ distinct techniques to generate the same conclusion: that identification of the counterfactual outcome requires analysts to control for the confounding variable V . Summarizing, we have that:

$$P_{fix}(Y(t) | V) = P_h(Y | \tilde{T} = t, V) = P_{do}(Y(t)|V) = P_e(Y | T = t, V). \quad (19)$$

If V were observed, we would be able to evaluate the expected value of the counterfactual outcome expectation, $E_h(Y | \tilde{T} = t)$, by integrating the observed expectation $E_e(Y | T =$

³¹See Heckman and Pinto (2015) for further discussion of the connection between empirical and hypothetical models.

t, V) over the support of V . The econometric literature provides a rich menu of strategies to control for the confounding variable V . We discuss part of this menu in the next section.

4 Econometric Approaches to Identification of Counterfactuals in the Generalized Roy Model

The Generalized Roy model is a framework for exploring the large toolkit of the econometric approaches for identifying counterfactuals. We compare what is possible in the econometric approach with what can be obtained using the non-econometric paradigms. We describe several of these approaches here. We develop this discussion further in subsequent sections of this paper.

Equation (19) states that the identification of causal effects in the Generalized Roy model hinges on controlling for the unobserved confounding variables V . A popular approach to doing so uses instrumental variables that are independent of V . It controls for V by shifting T without affecting the distribution of V . However, the IV approach using Z as an instrument does not identify interesting counterfactuals without additional assumptions.

For a simple example, consider a linear model in which the structural treatment equation is $T = \alpha_0 + \alpha_1 Z + \alpha_2 V + \epsilon_T$, and the outcome function is $Y = \beta_0 + \beta_1 T + \beta_2 V + \epsilon_Y$, where $\alpha_0, \alpha_1, \alpha_2, \beta_0, \beta_1, \beta_2$ are scalar parameters. In this model, the causal effect of T on Y is given by β_1 and is identified by the covariance ratio $cov(Y, Z)/cov(T, Z)$. This parameter can be estimated by a Two-Stage Least Squares (2SLS) Regression. This tool has been available to economists since the 1950s.³²

However, the Generalized Roy model is not captured by this simple two-equation system. The causal effect, $Y(1) - Y(0)$ is, in general, a random variable and not a constant so that treating β_1 as a constant does not capture essential heterogeneity of treatment effects across

³²See Amemiya (1985); Hansen (2022); Theil (1953, 1958, 1971). Theil (1953) invented this method. The method is far more general and applies to nonlinear models as well.

agents.³³ The analogue to heterogeneous β_1 is stochastically dependent on V . There are numerous approaches to identifying its distribution. We start with the use of instrumental variables in the presence of heterogeneous treatment effects and then consider alternative approaches.

4.1 Instrumental Variables

Heckman and Vytlacil (1999, 2005) address the question of identifying the Roy model by assuming a separable choice equation. Their approach enables analysts to control for V and, in turn, identify counterfactual outcomes. Their local Instrumental Variable (LIV) approach considers a binary treatment $T \in \{0, 1\}$. Their *separability assumption* is motivated by economic choice theory and states that treatment is given by a latent threshold-crossing equation that includes instrument Z and the confounder V ; that is, $T = \mathbf{1}[\zeta(Z) \geq \phi(V)]$. Separability enables them to rewrite the choice equation as:

$$T = \mathbf{1}[P(Z) \geq U]; \quad P(Z) = P_e(T = 1 | Z), \quad (20)$$

where the probability of treatment selection $P(Z) = P_e(T = 1 | Z)$ is sometimes called the propensity score. The unobserved variable U is given by $U = F_{e,\phi(V)}(\phi(V))$ where $F_{e,\phi(V)}$ is the cdf of $\phi(V)$, which is monotone increasing by construction. Subscript “ e ” denotes that the probability distribution is constructed using the empirical model. Variable U has a uniform distribution if $\phi(V)$ is absolutely continuous; that is, $U \sim \text{unif}([0, 1])$. The structural approach uses unobservables. The Neyman-Rubin approach does not. The *do-calculus* uses them, but in a limited way. We show in Section 6 that it rules out exploiting the information used to obtain (20). This approach to unobservables precludes the use of methods that are fruitful in the econometric approach.

³³A heterogeneous treatment effect case would write $\beta_1 = (Y_1 - Y_0)$ and $\beta_0 = Y_0$.

The hypothetical and empirical models for the Generalized Roy model that include the unobserved variable U are displayed in Table 5. The LMC of T in the hypothetical Roy model of Table 5 implies that $Y \perp\!\!\!\perp T \mid (Z, \tilde{T}, U)$. The LMC of Z implies $Y \perp\!\!\!\perp Z \mid (U, \tilde{T})$. These two independence relationships imply, by contraction property D, that $Y \perp\!\!\!\perp T \mid (\tilde{T}, U)$. Following the same analysis of V as (19), $Y \perp\!\!\!\perp T \mid (\tilde{T}, U)$ implies that:

$$P_h(Y \mid \tilde{T} = t, U) = P_{fix}(Y(t) \mid U) = P_{do}(Y \mid T = t, U). \quad (21)$$

Otherwise stated, controlling for U enables analysts to identify counterfactual outcomes in the same fashion that controlling for V does. Variable U is called a *balancing score* for V . This means that U is a surjective function of V that preserves the independence relationship $Y \perp\!\!\!\perp T \mid (\tilde{T}, V) \Rightarrow Y \perp\!\!\!\perp T \mid (\tilde{T}, U)$.³⁴

Table 5: Binary Choice Roy Model: Empirical and Hypothetical Causal Models		
	Empirical Model	Hypothetical Model
	<pre> graph TD V((V)) --> U((U)) V --> T[T] U --> T Z[Z] --> T T --> Y[Y] </pre>	<pre> graph TD V((V)) --> U((U)) V --> Y[Y] U --> T[T] Z[Z] --> T T --> Y Ttilde((T-tilde)) --> Y </pre>
	LMC	LMC
$V :$	$V \perp\!\!\!\perp Z$	$V \perp\!\!\!\perp (Z, \tilde{T})$
$Z :$	$Z \perp\!\!\!\perp (U, V)$	$Z \perp\!\!\!\perp (V, U, Y, \tilde{T})$
$U :$	$U \perp\!\!\!\perp Z \mid V$	$U \perp\!\!\!\perp (Y, Z, \tilde{T}) \mid V$
$T :$	$T \perp\!\!\!\perp V \mid (Z, U)$	$T \perp\!\!\!\perp (\tilde{T}, V, Y) \mid (Z, U)$
$Y :$	$Y \perp\!\!\!\perp (Z, U) \mid (T, V)$	$Y \perp\!\!\!\perp (Z, U, T) \mid (\tilde{T}, V)$
$\tilde{T} :$	(not defined for the model)	$\tilde{T} \perp\!\!\!\perp (T, V, U, Z)$

The Local Instrumental Variable (LIV) model of Heckman and Vytlacil (1999) can be used to identify probability distributions of counterfactual outcomes conditioned on U by taking the derivative of the observed outcome with respect to the propensity score. More generally, the counterfactual expectation $E_{fix}(g(Y(t)) \mid U = u)$ for any real-valued function

³⁴The balancing score was introduced by Rosenbaum and Rubin (1983).

$g : \mathbb{R} \rightarrow \mathbb{R}$ is identified if there is sufficient variation of propensity score $P(Z)$ around the value $u \in (0, 1)$.

Identification of $E_h(g(Y | \tilde{T} = t, U = u))$ comes from the derivative of the expectation $(-1)^{1-t} E_e(g(Y) \mathbf{1}[T = t] | P(Z))$ with respect to the propensity score at the value $P(Z) = u$. In particular, it can be shown that:

$$\begin{aligned} & E_h(Y | \tilde{T} = 1, U = u) - E_h(Y | \tilde{T} = 0, U = u) \\ & \equiv E_{fix}(Y(1) - Y(0) | U = u) = \frac{\partial E_e(Y | P(Z))}{\partial P(Z)} \Bigg|_{P(Z)=u} \end{aligned} \quad (22)$$

where *fix* refers to the distribution generated by fixing (which is the same as that generated by “doing”) and *e* refers to the sample distribution. Identification requires sufficient variation of the propensity score $P(Z)$ around $u \in [0, 1]$. If $P(Z)$ has full support, the average treatment effect can be evaluated by $ATE \equiv E_h(Y | \tilde{T} = 1) - E_h(Y | \tilde{T} = 0) = \int_0^1 (E_h(Y | T = 1, U = u) - E_h(Y | T = 0, U = u)) du$.

4.2 Stratification

A recurrent theme in this section is that identification of counterfactual outcomes hinges on controlling for the confounding variable V . The solution of the LIV model invokes separability assumption (20) which generates a balancing score U for V . According to (23), the nonparametric point-identification of the counterfactual outcomes conditioned on $U = u$ is obtained by differentiating the outcome with respect to the propensity score $P(Z)$ at value $u \in (0, 1)$.

Equation (22) assumes that the sample propensity score has enough variation around the value $u \in (0, 1)$. Consequently, the equation is not directly applicable to the case of discrete instruments. One approach to overcoming this limitation is to use the discrete counterpart of equation (22). Heckman and Vytlacil (2005) show that for any two values $z, z' \in \text{supp}(Z)$ such that $P(z') = u' > u = P(z)$ we have that:

$$\frac{E_e(Y | Z = z') - E_e(Y | Z = z)}{P_e(T = 1 | Z = z') - P_e(T = 1 | Z = z)} = \frac{\int_u^{u'} E_{fix}(Y(1) - Y(0) | U = u) du}{u' - u} \quad (23)$$

$$= E_{fix}(Y(1) - Y(0) | u \leq U \leq u').$$

Equation (23) states that difference of mean outcomes conditional on two instrumental values z, z' identifies the counterfactual outcome over an interval of U defined by the propensity scores $P(z)$ and $P(z')$. The equation evaluates a causal effect that depends on the values of the instrument. These effects are called Local Average Treatment Effects (LATE) by [Imbens and Angrist \(1994\)](#). LATE-type effects differ from causal effects such as ATE or TT, which do not depend on the IV values.³⁵

A consequence of (23) is that *ATE* can be identified if there are two instrumental variable values z_0, z_1 such that z_0 induces full treatment nonparticipation ($P(z_0) = 0$), and z_1 induces full treatment participation ($P(z_1) = 1$):

$$E_e(Y | Z = z_1) - E_e(Y | Z = z_0) = E_{fix}(Y(1) - Y(0) | 0 \leq U \leq 1)$$

$$= E_h(Y | \tilde{T} = 1) - E_h(Y | \tilde{T} = 0) = ATE.$$

This setup is equivalent to a randomized control trial with full compliance. [Mogstad and Torgovitsky \(2018\)](#) use functional form assumptions to extrapolate estimates over intervals of U to point estimates.

Another approach for controlling for V exploits the discrete nature of the instrument to generate an alternative balancing score. Let instrument Z take values in the discrete set $\text{supp}(Z) = \{z_1, \dots, z_N\}$ such that $P(z_1) < \dots < P(z_N)$.³⁶ Let $T(z) = \mathbf{1}[\zeta(z) \geq \phi(V)]$ be the counterfactual choice that would occur if Z were fixed at value $z \in \{z_1, \dots, z_N\}$. The *response vector* $\mathbf{S} = [T(z_1), \dots, T(z_N)]'$ is the random vector of potential choices across all Z -values.

³⁵[Heckman et al. \(2008\)](#) develop the relationship between LIV and LATE in depth.

³⁶The increasing ordering of propensity scores is assumed without loss of generality.

Response vector \mathbf{S} shares the same causal relationships of unobserved variable U in Table 5. By this we mean that \mathbf{S} is a function of V and that the choice T can be written as function of Z and \mathbf{S} :

$$T = \left[\mathbf{1}[Z = z_1], \dots, \mathbf{1}[Z = z_N] \right] \cdot \mathbf{S}.$$

Similar to U , the response vector \mathbf{S} is a balancing score for V . The independence relationship $Y \perp\!\!\!\perp T \mid (\tilde{T}, \mathbf{S})$ holds, which implies that $P_h(Y \mid \tilde{T} = t, \mathbf{S}) = P_e(Y \mid T = t, \mathbf{S})$. Heckman and Pinto (2018) show that the response vector \mathbf{S} controls for V by generating a special partition of its support that spans the support of V and renders choice T statistically independent of V within each cell of the partition. Each column of \mathbf{S} is a list of responses to different treatments for a person of a given V .

The values of \mathbf{S} are called response-types or strata.³⁷ The separability assumption eliminates some of potential response-types. An influential example is due to Imbens and Angrist (1994), who investigate the case of a binary instrument and a binary treatment. There are four possible response-types termed always-takers, compliers, never-takers and defiers. They invoke a monotonicity condition that is equivalent to the separability assumption Vytlacil (2002). The assumption eliminates the defiers and enables the identification of treatment effects for the compliers. See Heckman and Pinto (2018) and Buchinsky and Pinto (2021) for general results on identification.

4.3 The Matching Assumption

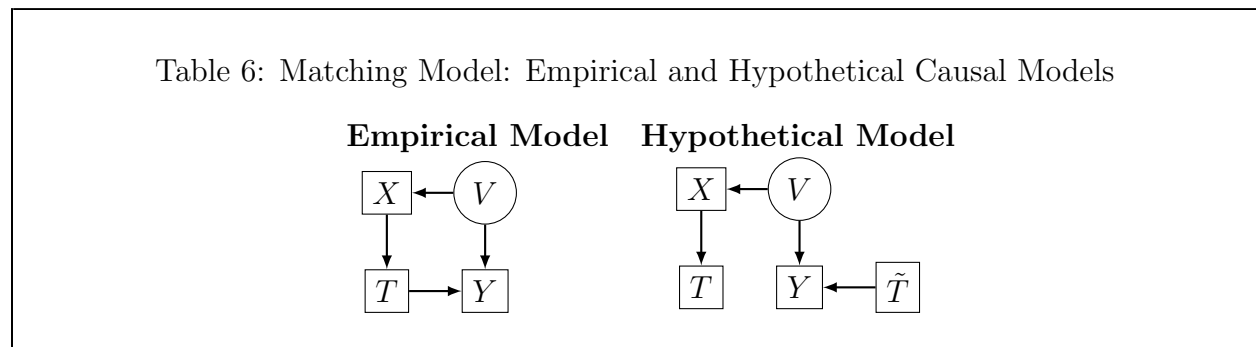
A popular method for identifying treatment effects assumes that a set of observed pre-treatment variables suffice to control for the confounding variable V . Otherwise stated, it assumes that the observed variable X is a balancing score for the confounding variable V . This assumption is called *Matching*.³⁸ Another (structural) way to state this is that X spans the space of V .

³⁷The concept was developed by Robins (1986) and embellished in Frangakis and Rubin (2002).

³⁸Heckman et al. (1998) investigate several estimation methods that invoke the matching assumption.

Table 6 presents the empirical and the hypothetical models that justify the matching assumption. The LMC of T in the hypothetical model implies that $Y \perp\!\!\!\perp T \mid (\tilde{T}, X)$. According to (17), we have that $P_h(Y \mid \tilde{T} = t, X) = P_{fix}(Y(t) \mid X) = P_e(Y \mid T = t, X)$ which means that the counterfactual outcome is identified by conditioning on X . Matching variables X are assumed not to be descendants of the hypothetical variable \tilde{T} . Thus, $P_h(X) = P_e(X)$ and the probability distribution of the counterfactual outcome is given by $P_{fix}(Y(t)) = \int (P_e(Y \mid T = t, X = x) dF_{e,X}(x))$.

The average causal effect of a binary treatment $T \in \{0, 1\}$ is evaluated by the weighted average of mean difference between the treated and not-treated participants that *match* on X , namely, $ATE = \int (E_e(Y \mid T = 1, X = x) - E_e(Y \mid T = 0, X = x)) dF_{e,X}(x)$.³⁹



The matching assumption replaces the *unobserved* variable U of the Generalized Roy model in Table 5 by the *observed* variable X . In practice, it assumes that potential bias generated by confounding variables can be ignored when controlling for observed pre-treatment variables. Under matching, the identification of treatment effects does not require an instrumental variable nor additional assumptions such as separability. This assumption enables us to solve the problem of selection bias induced by unobserved variables V by conditioning on the observed variables X .

³⁹Heckman et al. (1998) incorporated additive separability between observable and unobservable variables as well as exogeneity conditions that isolate outcomes and treatment participation into the matching framework. Additionally, they compare various types of estimation methods to show that kernel-based matching and propensity score matching have similar treatment of the variance of the resulting estimator.

The matching assumption is justified in the case of randomized controlled trials (RCTs). In this case, the matching variables X denote the pre-treatment variables conditioned on in the randomization protocol. In observational studies, the matching assumption is often rather strong. It assumes that the analyst observes enough information to make all the agent’s unobserved variables irrelevant (see Heckman, 2008b). Otherwise stated, matching assumes a symmetry in information between the economic agent and the econometrician.

There are several identification approaches that acknowledge the possibility of information asymmetries between the agent being studied and the econometrician: control function approaches, replacement functions or proxy variables. These methods often differ considerably in terms of assumptions and methodology. However, they all share the same identification principle: they use observed data to evaluate a proxy variable that plays the role of a matching variable.

4.4 Matching on Proxied Unobservables

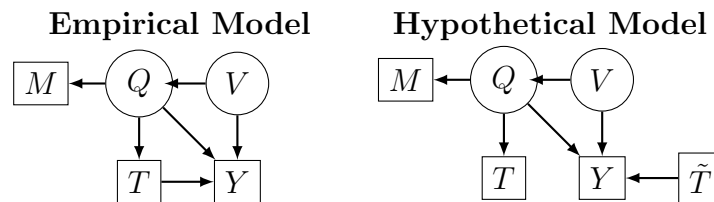
Matching on proxied unobservables is a version of matching that uses observed data to control for the confounding effects of V . Consider the modification of the Generalized Roy model in Table 7. The unobserved variable Q is a balancing score for the unobserved confounder V . The matching conditions of hypothetical model, $Y \perp\!\!\!\perp T \mid (\tilde{T}, Q)$, and its respective counterpart in the empirical model, $Y(t) \perp\!\!\!\perp T \mid Q$, hold. Variable Q has two additional properties: (1) it may cause outcome Y ; and (2) it may be measured with error by the observed variable M .

A common setup where Q arises is in the evaluation of college returns where T denotes college graduation, Y denotes earnings, and Q denotes unobserved abilities such as cognition or conscientiousness. These abilities are not directly observed but measured with error by an observed vector of variables M , such as psychological surveys or test scores. Formally, we write $M = f_M(Q, \epsilon_M)$. The identification strategy is to exploit the structural function

$M = f_M(Q, \epsilon_M)$ to evaluate Q , which, in turn, allows us to control for V and identify causal effects.

Matching on proxied unobservables has long been used in the economics of education (see, e.g., the essays in [Goldberger and Duncan, 1973](#) and [Goldberger, 1972](#)). The method is called the latent variable approach by [Heckman and Robb \(1985a\)](#). This literature offers several possibilities for estimating Q ([Aakvik et al., 1999, 2005](#); [Carneiro et al., 2003](#); [Cunha et al., 2005](#)). [Olley and Pakes \(1996\)](#) is an application of this method. A common parametric approach extracts factors from psychological measurements to extract Q as a latent factor. Nonparametric factor analysis is developed in [Cunha et al. \(2010\)](#) and [Schennach \(2020\)](#). It is also possible to condition nonparametrically on Q without knowing the functional form of f_M .

Table 7: Matching on Proxied Unobservables: Empirical and Hypothetical Causal Models



4.5 Control Functions

The control function principle specifies the dependence of the relationship between observables and unobservables in a nontrivial fashion. The principle was introduced in [Heckman and Robb \(1985a,b\)](#) building on earlier work by [Telser \(1964\)](#) and later popularized by [Blundell and Powell \(2003\)](#). It was also applied in [Carneiro et al. \(2003\)](#) and [Cunha et al. \(2005\)](#). Heckman’s sample selection correction ([1979](#)) is a control function.

We illustrate the control function principle using a version of the Generalized Roy model where V is a scalar random variable and the binary choice T is given by the *separable* equation $T = \mathbf{1}[\mu(Z) \geq V]$. Let $J = f_J(T, V, \epsilon_J)$ represents unobserved skills caused by the treatment

T and the unobserved confounding variable V . In addition, let the outcome equation be *additive* in K , that is to say that the outcome Y can be written as $Y = f_Y(T, \epsilon_Y) + \psi(J)$, The model is displayed as a DAG in Table 8. The LMC of Y in the hypothetical model implies that $Y \perp\!\!\!\perp T \mid (\tilde{T}, K)$. This means that K is a matching variable. The control function approach seeks to control for variable V by estimating the function $\psi(J)$ of the outcome equation.

Heckman and Vytlacil (2007a,b) use the assumption of separability of observables and unobservables in the choice equation and the outcome assumption of additivity to evaluate $\psi(J)$ as a function of the propensity score $P(Z)$. Similar to the LIV Model, we can use the CDF transformation to write the choice equation as $T = \mathbf{1}[P(Z) \geq F_V(V)]$, where $F_V(V) \sim \text{unif}([0, 1])$. Note that the expected value of the outcome conditional on $T = 1$ gives the *conditional* counterfactual mean:

$$E_e(Y \mid Z, T = 1) = E_{fix}h(Y(1) \mid Z, T = 1) = E_h(Y \mid \tilde{T} = 1, Z, T = 1),$$

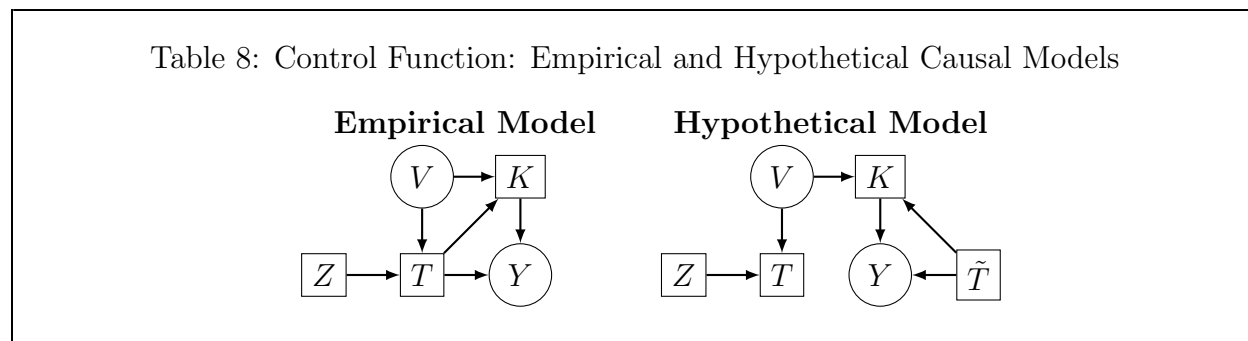
where the first term is observed, the second term uses fixing and the last one uses the hypothetical model. Under separability and outcome additivity, we can express $E_h(Y(1) \mid \tilde{T} = 1, Z, T = 1)$ as:

$$\begin{aligned} E_h(Y \mid \tilde{T} = 1, Z = z, T = 1) &= E_h(f_Y(\tilde{T}, \epsilon_Y) \mid \tilde{T} = 1) + E_h(\psi(J) \mid \tilde{T} = 1, Z = z, T = 1), \\ &= E_h(f_Y(1, \epsilon_Y)) + E_h(\psi(f_K(1, V, \epsilon_K)) \mid Z = z, T = 1), \\ &\quad \left(\text{setting } E_h(f_Y(1, \epsilon_Y)) = \alpha_1 \right) \\ &= \alpha_1 + E_h\left(\psi(f_K(1, V, \epsilon_K)) \mid P(z) > F_V(V)\right), \\ &= \alpha_1 + E_e\left(\psi(f_K(1, V, \epsilon_K)) \mid P(z) > F_V(V)\right). \\ \therefore E_h(Y \mid \tilde{T} = 1, Z, T = 1) &= \alpha_1 + \underbrace{f_1(P(Z))}_{\text{control function}}, \text{ where } f_1(P(Z)) \\ &= E_h(\psi(f_K(1, V, \epsilon_J)) \mid Z, T = 1). \end{aligned}$$

where the first equality uses the additivity assumption, the second uses the fact the \tilde{T} is an external variable, the third uses the separability assumption, the fourth switches the

hypothetical model into the empirical model as V , ϵ_K , Z are non-descendants of \tilde{T} . The last equation gives the expectation $E_h(Y | \tilde{T} = 1, Z, T = 1)$ as a function of the propensity score $P(Z)$. Control function $f_1(P(Z))$ can be estimated from observed data and the expected value of the counterfactual outcome can be evaluated as

$$E_h(Y(1)) = \int_0^1 \alpha_1 + f_1(p) dF_{P(Z)}(p).$$



4.6 Panel data Analysis and Other Approaches

A commonly used panel data method is **difference-in-differences** as discussed in Heckman and Robb (1985a), Blundell et al. (1998), Heckman et al. (1999), and Bertrand et al. (2004). All of the estimators previously discussed can be adapted to a panel data setting. Heckman et al. (1998) introduce difference-in-differences matching estimators to eliminate the bias in estimating treatment effects. Abadie (2005) extends this work. Separability between errors and observables is a common feature of the panel data approach in its standard application. Altonji and Matzkin (2005) and (Matzkin, 1993) present analyses of nonseparable panel data methods. Regression discontinuity estimators, which are versions of IV estimators, are discussed by Heckman and Vytlacil (2007b).

Table 9 summarizes some of the main identification approaches for the Generalized Roy model discussed here. The table barely scratches the surface, but gives a sense of the broad menu in the econometric approach. The essays in the *Handbooks of Econometrics* (Heckman and Leamer, 2001, 2007) give a range of other estimation approaches.

Table 9: Some Alternative Approaches that Identify Treatment Effects by Controlling for V

$$Y \perp\!\!\!\perp T \mid (\tilde{T}, X, V), \quad T \in \{0, 1\}$$

$$E_h(Y \mid \tilde{T} = t, X = x) = \int E_e(Y \mid T = t, X = x, V = v) dF_{e,V \mid X=x}(v)$$

	Method Assumes	Need Instrument (Z)?	Identify Distribution of V ?
Matching ^a	V, X known	No	Yes (V observed)
Control Functions ^b	V estimated, X, Z known (continuous T); Bounds on quantiles of V estimated (discrete case)	Yes (except cases where functional forms secure identification)	Yes (over support)
Factor Method ^c	Distribution of V estimated from additional measurements of V (M)	No	Yes (with auxiliary measurements over support)
IV: LATE, LIV ^d	Z, X known	Yes	Estimate intervals of quantiles of V (Heckman and Vytlacil, 1999, 2005) and conditions on them; LIV shrinks interval of quantiles of V to a point using continuous instruments and conditions on them
Stratification ^e	Z, X known	Instruments give restrictions on strata (balancing scores for V)	Identify distribution of strata which places interval bounds on V and conditions on them
Longitudinal Data Methods ^f	Variety of assumptions	Covariance restrictions	Yes; in long panels can identify V
Mixing Distributions ^g	$V \perp\!\!\!\perp X$	No (intervals of V)	Yes (Mixtures)

^a Heckman et al. (1998); Rosenbaum and Rubin (1983); ^b Blundell and Powell (2003); Heckman and Robb (1985a,b); ^d See review in Heckman and Vytlacil (2007a); ^e Frangakis and Rubin (2002); Heckman and Pinto (2018); ^f Abbring and Heckman (2007); Heckman and Robb (1985a); ^g Cameron and Heckman (1998); Heckman and Singer (1984); Prakasa Rao (1992)

5 The Neyman-Rubin (NR) Causal Model

The Neyman-Rubin causal approach uses the language and framework of experimental design developed by [Neyman \(1923\)](#), [Fisher \(1935\)](#), and [Cox \(1958\)](#) and popularized by [Holland \(1986\)](#). It ignores essential aspects of the econometric approach to causality and conflates distinct concepts (e.g., SUTVA).⁴⁰ It does not define hypothetical models nor does it employ structural equations to characterize causal models. It focuses on units of analysis instead of system of equations.

In this approach, causal models are characterized by statistical independence relationships among counterfactual counterparts of observed variables, never precisely justified. In place of the thought experiments that characterize the econometric approach, it uses randomized controlled experiments as the foundational paradigm. The contrast between thought experiments and randomized control experiments is central to understanding the differences in the approaches.

The NR approach lacks the clarity of interpretation offered by causal models described by structural equations. It is often difficult to map the independence relationships of a NR model into the causal relationships produced by economic theory. In particular, NR makes it difficult to use economic theory to assess the credibility of assumptions about the underlying structural equations that ensure the identification of causal effects or to interpret economic data using economic models.

Another drawback is that the NR framework lacks the fundamental tools of econometric causal analysis. It does not explicitly model unobserved variables in structural models. This feature substantially limits the use of standard econometric tools. It rules out (or makes cumbersome) several fruitful econometric strategies such as balancing bias within models using compensating variations of arguments of structural functions to keep agents at the

⁴⁰[Heckman \(2008a\)](#) explains that: SUTVA - Stable Unit Treatment Value Assumption - conflates two two distinct concepts regarding functional autonomy (structural invariance) and no interactions among agents.

same levels of well being,⁴¹ and cross-equation restrictions on both observable and unobservable model components,⁴² or functional form restrictions. In practice, the set of tractable identification strategies that employ the NR framework is limited to a few possibilities: randomized trials, matching, IV and its many surrogates and differences-in-differences.⁴³ This section illustrates the drawbacks of NR in analyzing core policy questions or in synthesizing and interpreting evidence.

5.1 The Generalized Roy Model under NR

The NR framework focuses on the unit of analysis $i \in \mathcal{I}$ which usually represents an economic agent or entity. The framework describes part of the Generalized Roy model of Table 3 using two counterfactuals: $T_i(z)$ is the potential treatment when the instrument Z is externally set to value $z \in \text{supp}(Z)$; and $Y_i(t, z)$ is the potential outcome of agent i when Z is set to value $z \in \text{supp}(Z)$ and choice T is set to $t \in \text{supp}(T)$. Properly formulated, potential outcomes are the outputs of structural equations. NR does not explicitly characterize the treatment choice equation. It prides itself on being nonparametric, although some proponents claim that assuming linearity is an innocuous assumption, even when models are fundamentally nonlinear.⁴⁴

The NR framework characterises the Generalized Roy model (4)–(7) by three assumptions:

1. An exclusion restriction states that $Y_i(t, z) = Y_i(t, z')$, for all $z, z' \in \text{supp}(Z)$, $t \in \text{supp}(T)$ and all $i \in \mathcal{I}$.
2. IV relevance: Z is not statistically independent of T , that is $Z \not\perp T$.
3. Exogeneity condition: $Z \perp\!\!\!\perp (Y(t), T(z))$ for all $(z, t) \in \text{supp}(Z) \times \text{supp}(T)$.

⁴¹See e.g., Ekeland et al. (2004); Rosen (1986).

⁴²See, e.g., Hansen and Sargent (1982).

⁴³See Imbens and Rubin, 2015.

⁴⁴Angrist and Pischke (2009). Ekeland et al. (2004) show that nonlinearity is intrinsic to hedonic models and that linearizing it produces identification problems.

The exclusion restriction means that Z does not directly cause Y . Thus, we can express the counterfactual outcome as $Y_i(t)$ instead of $Y_i(t, z)$. IV relevance means that T is caused by Z . The exogeneity condition of the NR framework can be traced back to the independence relationship between Z and V of the Generalized Roy model (4)–(7). In the NR framework, the exogeneity condition is an assumption. In the Generalized Roy model, the exogeneity condition is a consequence of the causal relation among model variables. Namely, that the Z and V are external variables. LMC (8) implies that $Z \perp\!\!\!\perp V$, which, in turn, generates the exogeneity condition.

The identification of counterfactual outcomes requires additional assumptions. A popular assumption securing identification is the monotonicity condition (24) of Imbens and Angrist (1994). It states that a change in an instrument induces agents to change their treatment choice in the same direction. Notationally, for any $z, z' \in \text{supp}(Z)$, it says that:

$$T_i(z) \geq T_i(z') \quad \forall i \in \mathcal{I} \quad \text{or} \quad T_i(z) \leq T_i(z') \quad \forall i \in \mathcal{I}. \quad (24)$$

Vytlacil (2002) shows that the monotonicity condition (24) is equivalent to the separability assumption $T = \mathbf{1}[\zeta(Z) \geq \phi(V)]$. Otherwise stated, the NR counterpart for the Generalized Roy model separability assumption is the monotonicity condition. Each condition enables the identification of causal effects of T on Y in its respective framework. At this level, the IV models in the two frameworks are equivalent.

Model equivalence does not, however, imply that they offer the same analytical capacities. In particular, the Generalized Roy model (4)–(7) explicitly displays the unobserved confounding variable V , while NR does not. This feature enables analysts to further investigate the model and use other approaches for controlling for it. Section 4 shows that the identification of counterfactual outcomes hinges on the analysts’s ability to control for the unobserved confounding variable V . Heckman and Vytlacil (2005) use the fact that U is a balancing score for V to define and identify a new parameter called the marginal treatment

Table 10: Some Causal Parameters as Weighted Average the MTE

Causal Parameters	MTE Representation	Weights
$ATE = E_h(Y(1) - Y(0))$	$= \int_0^1 MTE(p)W^{ATE}(p)dp$	$W^{ATE}(p) = 1$
$TT = E_h(Y(1) - Y(0) T = 1)$	$= \int_0^1 MTE(p)W^{TT}(p)dp$	$W^{TT}(p) = \frac{1 - F_{e,P}(p)}{\int_0^1 (1 - F_{e,P}(t))dt}$
$TUT = E_h(Y(1) - Y(0) T = t_0)$	$= \int_0^1 \Delta^{MTE}(p)W^{TUT}(p)dp$	$W^{TUT}(p) = \frac{F_{e,P}(p)}{\int_0^1 (1 - F_{e,P}(t))dt}$
$TSLS = \frac{Cov(Y, Z)}{Cov(T, Z)}$	$= \int_0^1 MTE(p)W^{TSLS}(p)dp$	$W^{TSLS}(p) = \frac{\int_0^1 (t - E_e(P))dF_{e,P}(t)}{\int_0^1 (t - E_{e,P}(t))^2 dF_{e,P}(t)}$
$LATE = \frac{E_e(Y Z = z_1) - E_e(Y Z = z_0)}{P_e(z_1) - P_e(z_0)}$	$= \int_{P(z_0)}^{P(z_1)} MTE(p)W^{LATE}(p)dp$	$W^{LATE}(p) = \frac{1}{P_e(z_1) - P_e(z_0)}$

Source: [Heckman and Vytlacil \(2005\)](#).

effect (MTE):

$$MTE(u) = E_h(Y | \tilde{T} = 1, U = u) - E_h(Y | \tilde{T} = 0, U = u) = E_{fix}(Y(1) - Y(0) | U = u).$$

The MTE plays a primary role in generating a range of causal effects commonly sought in policy evaluations. A few of these causal parameters are presented in Table 10.

The analytical gain generated by switching from the NR framework to a structural equation framework is substantial. The use of structural equations facilitates a richer analysis and a deeper investigation of the properties of the Generalized Roy model. Such analyses cannot be achieved in the NR framework because it does not include unobserved variables, nor does it employ structural equations. This analytical deficiency of the NR framework limits the researcher's ability to extend causal analysis of the Generalized Roy model and other economic models.

The parsimonious machinery of the NR framework is often misunderstood as endowing the Generalized Roy model with a greater level of generality. This impression is misleading as the IV model featured in the NR framework is equivalent to the Generalized Roy

model described by equations (4)–(7) and its monotonicity criteria is equivalent to a separability condition. Its apparent simplicity is due to its lack of explicit statements about its assumptions.

5.2 The Matching Model in NR

A common identification approach in NR is a *matching* assumption on observed variables. It states that the treatment choice T is independent of counterfactual outcomes $Y(t)$ when conditioning on observed pre-treatment variables X , that is, $Y(t) \perp\!\!\!\perp T \mid X$.⁴⁵ Intuitively, the assumption states that pre-treatment variables X are sufficiently rich to account for all the unobserved variables that jointly influence treatment choice T and outcome Y . The assumption can be easily criticized as often being overly optimistic for the case of observational studies (Heckman, 2008b; Heckman and Navarro, 2004).

It is natural to assume that increasing the number of variables on which matching is based decreases the bias generated by unobserved confounders. This statement is known to be false.⁴⁶ However it is rather difficult to investigate the truth of this claim using the NR framework. The causal model of Table 11 clarifies this point.

Table 11: Hypothetical Matching Model

Causal Model	DAG	Independence Relationships
$V = f_V(\epsilon_V)$ $J = f_J(\epsilon_J)$ $W = f_W(\epsilon_W)$ $V = f_V(\epsilon_V)$ $T = f_T(V, W, \epsilon_T)$ $R = f_R(T, V, \epsilon_R)$ $U = f_U(K, \epsilon_U)$ $X = f_X(W, J, \epsilon_X)$ $Y = f_Y(T, K, U, J, \epsilon_Y)$	<pre> graph TD V((V)) --> K(K) K --> U((U)) V --> T(T) T --> Y(Y) W((W)) --> X(X) J((J)) --> X X --> T U --> Y J --> Y </pre>	$Y(t) \perp\!\!\!\perp T \mid R$ $Y(t) \not\perp\!\!\!\perp T \mid X$ $Y(t) \not\perp\!\!\!\perp T \mid (X, R)$

⁴⁵In the language of Pearl (2009b), X *d-separates* Y and T .

⁴⁶See, for instance, Greenland et al. (1999); Heckman and Navarro (2004); Pearl (2009c).

The causal model Table 11 consists of four observed variables: the treatment T , the outcome Y , a pre-treatment variable X and a post-treatment variable R . The model also contains four unobserved variables V, U, W, J . The causal relationships among the observed and unobserved variables renders $Y(t) \perp\!\!\!\perp T \mid R$ even though $Y(t) \not\perp\!\!\!\perp T \mid X$. The independence relationship that characterises the matching assumption holds for post-treatment variables, but not for the pre-treatment variable. Moreover, adding the pre-program variable X to the conditioning set of $Y(t) \perp\!\!\!\perp T \mid K$ prevents identification because $Y(t) \not\perp\!\!\!\perp T \mid (X, K)$.

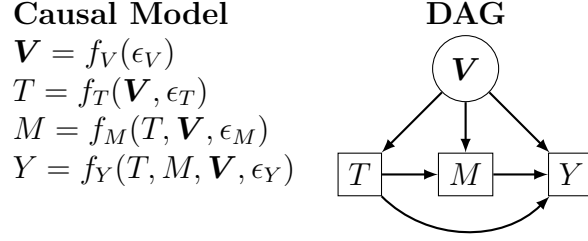
The causal model of Table 11 exemplifies the difficulty of performing causal investigations within the NR framework. The unusual properties of the model stem from the particular causal relationships among its observed and unobserved variables. This model is not easily analyzed within the NR framework which it lacks unobserved variables and suppresses the structural equations that clearly describe the causal relationships among variables.

5.3 Mediation Models under NR: An example

Mediation models originate in the literature on path analysis and simultaneous equations.⁴⁷ They trace the impacts of interventions on outcomes through their multiple channels of operation. Identifying the causal models generated by NR assumptions is often a daunting task. The economic content of these assumptions is often far from clear. We examine several mediation models to illustrate this point and show the power of the econometric approach compared to an approach based on NR principles. Table 12 uses the econometric approach to present a general mediation model in which a treatment T causes a mediator M and an outcome Y that is caused by both T and M . \mathbf{V} denotes a random vector that plays the role of the unobserved confounder causing T , M and Y . The counterfactual mediator when the treatment is fixed at $t \in \text{supp}(T)$ is $M(t) = f_M(t, \mathbf{V}, \epsilon_M)$. The counterfactual outcome when the treatment is fixed at t and the mediator is fixed at $m \in \{0, 1\}$ is $Y(t, m) = f_Y(t, m, \mathbf{V}, \epsilon_Y)$. The counterfactual outcome when we fix only T at t is $Y(t) = f_Y(t, M(t), \mathbf{V}, \epsilon_Y)$.

⁴⁷See Bollen (1989); Klein and Goldberger (1955); Wright (1921, 1934).

Table 12: Mediation Model with Confounding Variable



The goal of mediation models is to decompose the total effect of T on Y into an indirect effect that includes the effect of T on M and M on Y and a direct effect not mediated by M . To facilitate the discussion, let T and M denote binary variables taking values in $\{0, 1\}$. The average (total) effect of T on Y is $E_{fix}(Y(1) - Y(0))$.⁴⁸ We can also define the average direct effect of T on Y as $E_{fix}(Y(1, M) - Y(0, M)) = \sum_{m=0}^1 E_{fix}(Y(1, m) - Y(0, m))P_e(M = m)$ and the average indirect effect as $E_{fix}(Y(T, 0) - Y(T, 1)) = \sum_{t=0}^1 E_{fix}(Y(t, 1) - Y(t, 0))P_e(T = t)$.⁴⁹ Table 13 displays three hypothetical models suitable for examining the total, direct and indirect effects. The first DAG corresponds to the total effect. The hypothetical variable \tilde{T} replaces the T -input of both the mediator M and the outcome Y equations. The second DAG corresponds to the indirect effect only and the hypothetical variable replaces only the T -input of the mediator equation. The last DAG corresponds to the direct effect only where the hypothetical variable \tilde{T} replaces only the T -input of outcome equation.

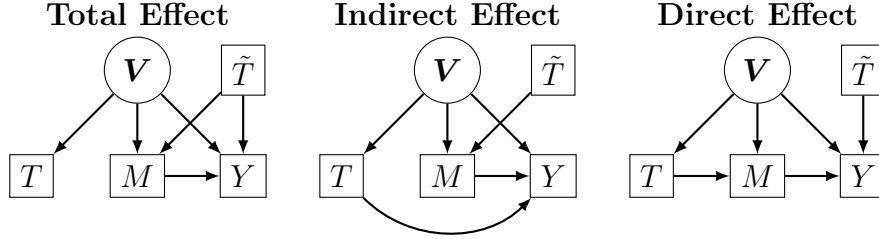
⁴⁸This is the same as $E_h(Y(1) - Y(0))$.

⁴⁹Alternatively, we can then define the direct effect and indirect effects for a given t by (25) and (26) respectively.

$$DE(t) = E_{fix}(Y(1, M(t)) - Y(0, M(t))) = \int E_{fix}(Y(1, m) - Y(0, m))dF_{M(t)}(m) \quad (25)$$

$$IE(t) = E_{fix}(Y(t, M(0)) - Y(t, M(1))) = \int E_{fix}(Y(t, m))dF_{M(1)}(m) - \int E_{fix}(Y(t, m))dF_{M(0)}(m). \quad (26)$$

Table 13: Hypothetical Models for the Mediation Model: Total, Direct and Indirect Effects



The presence of confounding variable V prevents the identification of the counterfactual means $E_{fix}(M(t))$ and $E_{fix}(Y(t, m))$. A solution to this identification problem using NR is the Sequential Ignorability (SI):⁵⁰

$$(Y(t', m), M(t)) \perp\!\!\!\perp T, \quad (27)$$

$$Y(t', m) \perp\!\!\!\perp M(t) \mid T, \quad (28)$$

for any $t, t' \in \text{supp}(T)$ and $m \in \text{supp}(M)$. SI (27)–(28) enables analysts to identify counterfactual means by statistical conditioning $E_e(M(t)) = E_{fix}(M \mid T = t)$ and $E_{fix}(Y(t, m)) = E_e(Y \mid T = t, M = m)$.

SI assumptions (27)–(28) can be understood as an application of the matching condition to mediation models. Assumption (27) states that the choice T is exogenous with respect to the outcome and mediator counterfactuals. The assumption would be justified if T were randomly assigned by a RCT experiment.

The interpretation of assumption (28) is less straightforward. It states that the counterfactual mediator $M(t)$ is independent of the counterfactual outcome $Y(t, m)$ when conditioned on T . The assumption cannot be directly tested even in randomized experiments that randomize T (Imai et al., 2010). SI assumptions (27)–(28) are much more easily interpreted using structural equations. The assumptions rule out any confounding variable V , generating the model in Table 14.

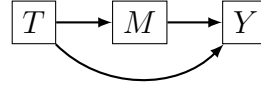
⁵⁰See Imai et al. (2011, 2010) for the properties of these assumptions. Levin and Robbins (1983) uses such assumptions in his g-computation algorithm.

Table 14: Mediation Model with No Confounding Variables

Causal Model

$$\begin{aligned} T &= f_T(\epsilon_T) \\ M &= f_M(T, \epsilon_M) \\ Y &= f_Y(T, M, \epsilon_Y) \end{aligned}$$

DAG



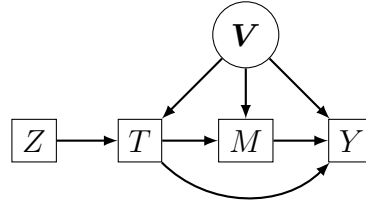
In light of a structural analysis, it can be seen that SI assumptions (27)–(28) are rather strong. They can be weakened if instrumental variables are available as depicted in Table 15. We use the model to exemplify a case in which NR assumptions are logically possible but generate a causal model that is difficult to justify using any plausible economic argument. The structural model enables the analyst to interpret the statistical assumptions using behavioral theory.

Table 15: Mediation Model with Instrumental Variables

Causal Model

$$\begin{aligned} \mathbf{V} &= f_V(\epsilon_V) \\ Z &= f_Z(\epsilon_Z) \\ T &= f_T(Z, \mathbf{V}, \epsilon_T) \\ M &= f_M(T, \mathbf{V}, \epsilon_M) \\ Y &= f_Y(T, M, \mathbf{V}, \epsilon_Y) \end{aligned}$$

DAG



The mediation model with IV has four counterfactuals, $T(z)$, $M(t)$, $Y(t)$, $Y(t, m)$ previously defined. In the language of NR, the model would be characterized by IV exogeneity condition $Z \perp\!\!\!\perp (T(z), M(t), Y(t), Y(t, m))$. The condition holds due to the independence of Z and \mathbf{V} .⁵¹ Suppressing Y generates an IV model where M plays the role of the outcome.

To dig more deeply, we investigate the case of a binary instrument $Z \in \{0, 1\}$. The response vector $\mathbf{S}_i = [T_i(0), T_i(1)]'$ denotes the vector of treatment choices that agent i would take if it were assigned to each of the instrumental values. Section 4 shows that, given \mathbf{S} , the treatment choice T depends only on the instrument Z . The exogeneity condition

⁵¹Note that if we were to suppress M from the DAG of Table 15, we would obtain the empirical model of Table 4.

states that Z is independent of the counterfactual outcome $Y(t)$. Thus

$$T \perp\!\!\!\perp Y(t) \mid \mathbf{S}. \quad (29)$$

\mathbf{S} is a balancing score for \mathbf{V} .

Yamamoto (2014) uses the language of NR to identify mediation effects using instrumental variables. His solution merges SI (27)-(28) with the matching property of the response vector \mathbf{S} in (29). He advocates an assumption that he terms the *local average causal mediation effects (LACME) assumption*:

$$(Y(t, m), M(t')) \perp\!\!\!\perp T \mid (\mathbf{S} = [0, 1]'), \quad (30)$$

$$Y(t, m) \perp\!\!\!\perp M(t') \mid (T, \mathbf{S} = [0, 1]'). \quad (31)$$

LACME (30)–(31) adds the the response vector \mathbf{S} as an additional conditioning variable to the SI independence relationships in (27)-(28). Assumption (30) is a simple extension of the matching property of \mathbf{S} from the IV model of Table 14 to the mediator model of Table 15. Under monotonicity (24), the LACME assumption identifies the direct and indirect mediation effects for compliers.

It is easy to interpret LACME in terms of NR assumptions: assumptions(30)–(31) are a weaker version of SI (27)-(28) that incorporates the LATE analysis of Imbens and Angrist (1994). On the other hand, it is difficult to gauge how the LACME assumptions fit into the mediation model of Table 12. It is even harder to interpret the causal content of these assumptions.

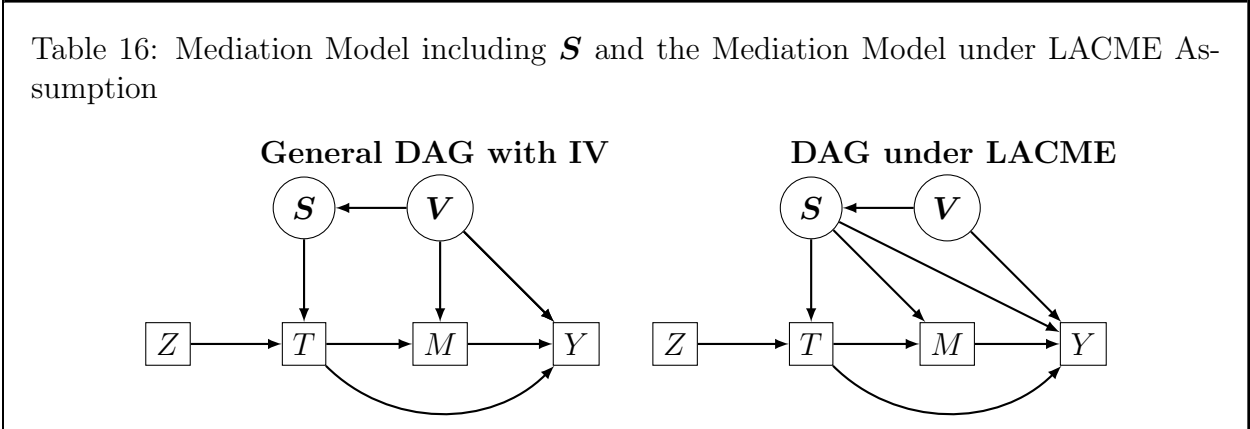
Table 16 presents two DAGs that use the structural approach to clarify the causal content of LACME. The first DAG places the unobserved response vector \mathbf{S} into the mediation model of Table 12. The response vector \mathbf{S} plays the role of a balancing score for \mathbf{V} only for choice T .⁵² The addition of the response vector does not result in any loss of generality. The second DAG displays the mediation model under LACME. According to assumption (31), the response vector \mathbf{S} plays the role of a balancing score for T and M . In addition, LACME

⁵²This property is based on the discreteness of the instrument.

prevents \mathbf{V} from jointly causing M, Y and implies that \mathbf{S} directly causes M, Y . It is hard to translate LACME into credible economic causal relationships.

$\mathbf{S} = [T(0), T(1)]'$ is expressed as a function of the confounding variable \mathbf{V} because $T(z)$ is a function of \mathbf{V} . Note that the choice T is expressed as a function of \mathbf{S} and Z because $T = [\mathbf{1}[Z = 0], \mathbf{1}[Z = 1]]\mathbf{S}$. The response vector $\mathbf{S} = [T(0), T(1)]'$ is expressed as a function of the confounding variable \mathbf{V} because $T(z)$ is a function of \mathbf{V} . The resulting DAG does not include more information than the original model of Table 12 because \mathbf{S} is unobserved.

The second DAG displays the mediation model under LACME. From assumption (31), the response vector \mathbf{S} plays the role of a matching variable for the causal effect of M on Y . It plays the role of a balancing score for \mathbf{V} for T, M , and Y . This assumption prevents \mathbf{V} from jointly causing M, Y and implies that \mathbf{S} directly causes M, Y . It is hard to produce interpretable economic models that justify \mathbf{S} as a cause of M or Y . LACME is an unmotivated statistical assumption devoid of economic content typical of analyses within the NR framework.



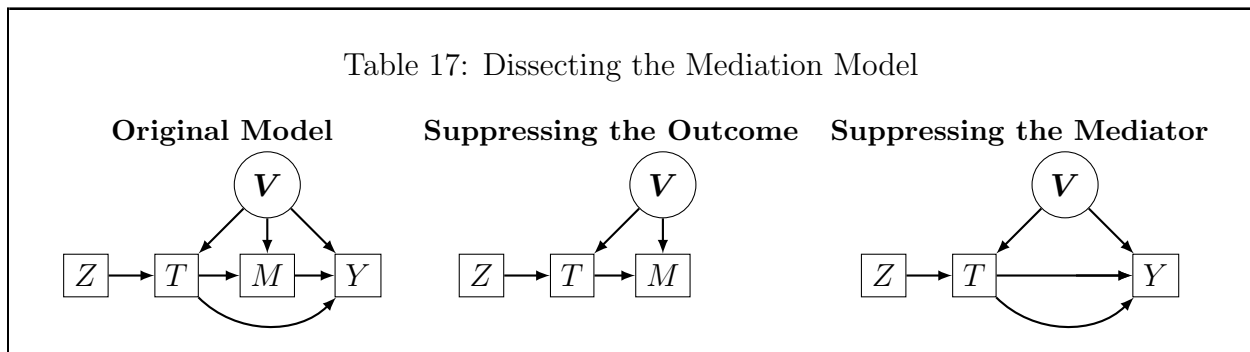
Using Structural Equations to Identify the Mediation Model with IV

Dippel, Gold, Heblich, and Pinto (2020) study the identification of causal effects for the mediation model with an instrumental variable. Their analysis illustrates the gain in clarity

and interpretability when a causal model is expressed by structural equations instead of NR statistical independence relationships.

A typical empirical setting of an IV model consist of one instrument and various outcomes. A mediation model with an instrument arises when treatment causes an intermediate outcome (the mediator), which in turn causes a final outcome. The DAG of this empirical model is presented in the first column of Table 17.

The second column of Table 17 presents the DAG generated by suppressing the final outcome. The resulting DAG is an IV model like that examined in Section 3. The causal effect of T on M can be identified by the methods discussed in Section 4. The third column of Table 17 suppresses the mediator M . The resulting model is also an IV model. This means that the *total effect* of T on Y can also be identified by the methods of Section 4. Unfortunately, the IV does *not* identify the causal effect of M on Y . Consequently, mediation analysis cannot be conducted without further assumptions.



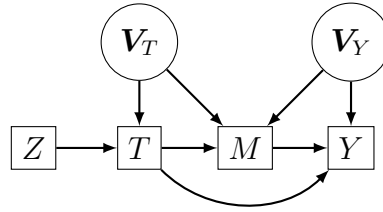
Dippel, Gold, Heblich, and Pinto (2020) address the question of whether it is possible to use an instrumental variable Z to nonparametrically identify the causal chain connecting T , M , Y while maintaining the endogeneity of the treatment T with respect to the mediator M and outcome Y . They show that the only solution to this problem is to assume the *partially confounded* mediation model of Table 18.

Table 18: Partially Confounded Model with Instrumental Variables

Causal Model

$$\begin{aligned} \mathbf{V}_T &= f_{V_T}(\epsilon_{V_T}) \\ \mathbf{V}_Y &= f_{V_Y}(\epsilon_{V_Y}) \\ Z &= f_Z(\epsilon_Z) \\ T &= f_T(Z, \mathbf{V}_T, \epsilon_T) \\ M &= f_M(T, \mathbf{V}_T, \mathbf{V}_Y, \epsilon_M) \\ Y &= f_Y(T, M, \mathbf{V}_Y, \epsilon_Y) \end{aligned}$$

DAG



The partially confounded assumption is that $\mathbf{V}_T \perp\!\!\!\perp \mathbf{V}_Y$. This is sometimes called a components of variance model. The assumption generates an additional exogeneity condition $(M(z), Y(m, t)) \perp\!\!\!\perp Z \mid (T = t)$ while maintaining the endogeneity of the treatment T with respect to M and Y . This means that Z is a valid instrument for identifying the causal effect of M on Y when conditioning on the treatment variable T . If the assumption holds, the causal effect of M on T can be evaluated by the methods of Section 4. [Dippel, Gold, Heblich, and Pinto \(2020\)](#) discuss the intuition, plausibility, and estimation of the partially confounded mediation model. They illustrate a range of examples where the partially confounding assumption may hold and where it does not.

6 The Do-Calculus and the Hypothetical Model

This section compares the *do*-calculus framework (DoC) of [Pearl \(2009b\)](#) with the Neyman-Rubin (NR) framework of [Holland \(1986\)](#) and [Imbens and Rubin \(2015\)](#) as well as the Hypothetical Model (HM) of [Heckman and Pinto \(2015\)](#).

The DoC was first presented in [Pearl \(1995\)](#). The method employs graph theory-based algorithms to identify the probability distribution of counterfactual variables in causal models represented by DAGs.⁵³ In contrast with NR, DoC uses autonomous (invariant) structural

⁵³For a recent book on the graphical approach to causality, see [Peters et al. \(2017\)](#). For related works on causal discovery, see [Glymour et al. \(2014\)](#), [Heckman and Pinto \(2015\)](#), [Hoyer et al. \(2009\)](#), and [Lopez-Paz et al. \(2017\)](#).

equations. The method clearly describes causal relationships among model variables. Its fundamental relationships are based on thought experiments. It is not incompletely formulated in a way that leads to problematic causal interpretations as in the NR approach.

DoC applies to any nonparametric, recursive system of structural equations. Similar to HM, DoC allows for unobserved variables. It can be applied to multiple equation causal models and a range of causal inquiries.

However, HM and the DoC differ greatly regarding counterfactual manipulations. To address the causal operation of fixing, the HM solution is based on a hypothetical model that formalizes thought experiments and places them on a sound probabilistic footing. Contrary to HM, DoC defines hypothetical models by making manipulations *within* the empirical model. It does not have a counterpart to \tilde{T} , the source of hypothetical variation in HM. DoC implements the notion of setting or fixing using a new set of rules that combine graphical analysis, independence relationships and probability equalities.

For instance, the DoC uses a DAG-based criteria called *d*-separation to check for conditional independence among variables. Its definition requires some DAG terminology. Let U be a path of arrows that connects variables T and Y in a DAG G regardless of the arrows' directions. A collider C in path U is a variable that has two arrows pointing at it (inverted fork). A variable V in the path U is said to block T and Y in the DAG G if it is not a collider (nor a descendant of a collider). T and Y are said to be *d*-separated by a set of variables V if V *d*-separates all paths from T to Y .

6.1 The Rules of DoC

As noted in Section 3.5, DoC uses the “back-door” criterion to verify matching conditions. In DoC terminology, the matching condition $Y(t) \perp\!\!\!\perp T|V$ of the Generalized Roy Model in Table 3 is expressed by the statement: “ V *d*-separates Y and T in the DAG $G_{\underline{T}}$,” where G be the original DAG of the Roy Model and $G_{\underline{T}}$ is a derived DAG which suppresses the arrows

departing from T . The “back-door” criterion holds for confounder V of the Roy Model. This implies the matching condition in which controlling for V renders the counterfactual outcome $Y(t)$ statistically independent of treatment T .

The core machinery of the DoC consists of three DAG-based rules. Additional notation is necessary to describe these rules. Let Y, K, Z, T denote disjoint variable sets in \mathcal{T} . In DoC notation, $T(Z)$ denotes the variables in T that do not directly or indirectly cause Z . “Do” deletes certain links in the original graph and assumes certain conditional independence relations. This is Pearl’s way to fix variables externally. DoC uses $G_{\bar{K}}$ for the derived DAG that deletes all causal arrows *arriving* at K in the original DAG G . $G_{\underline{T}}$ denotes the DAG that deletes all causal arrows emerging from T . In this notation, $G_{\bar{K},\underline{T}}$ stands for the derived DAG that suppresses all arrows *arriving* at K and emerging from T , while $G_{\bar{K},\overline{T(X)}}$ deletes all arrows arriving at K in addition to arrows arriving at $T(X)$, variables in T that are not ancestors of X . The DoC rules combine a graphical condition and a conditional independence relation that, when satisfied, imply a probability equality:

Table 19: The Three DoC Rules

Rule 1:

If $Y \perp\!\!\!\perp T \mid (K, Z)$ holds in $G_{\bar{K}}$, then $P(Y \mid do(K), T, Z) = P(Y \mid do(K), Z)$

Rule 2:

If $Y \perp\!\!\!\perp T \mid (K, Z)$ holds in $G_{\bar{K},\underline{T}}$, then $P(Y \mid do(K), do(T), Z) = P(Y \mid do(K), T, Z)$

Rule 3:

If $Y \perp\!\!\!\perp T \mid (K, Z)$ holds in $G_{\bar{K},\overline{T(Z)}}$, then $P(Y \mid do(K), do(T), Z) = P(Y \mid do(K), Z)$

The process of checking if a causal effect is identified requires iterative use of these rules. We present several examples of how to use these rules below.

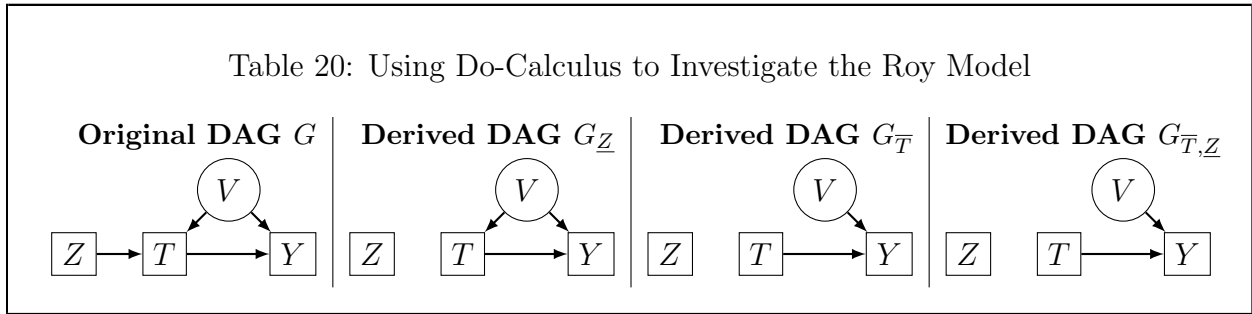
In computer science, DoC is said to be “complete.” This is different from the notion of completeness as defined in simultaneous equations theory discussed below in Section 7. The key DoC notion is that if a causal effect is identifiable, it can be identified by the iterative

application of some sequence of the three rules (Huang and Valtorta, 2006; Shpitser and Pearl, 2006).

A major limitation of do-calculus is that it only applies to non-parametric models that can be fully characterized by a DAG. Stated otherwise, the method does not account for assumptions about the functional forms of the structural equations or covariance restrictions. This limitation hinders the application of most of the popular econometric tools used in empirical economics such as cross equation restrictions, separability, additivity or monotonicity assumptions. For instance, the Generalized Roy model is not identified by DoC because it requires assumptions such as separability. The same is true of the IV model. Separability cannot be characterized by conditional independence assumptions generated by a DAG. We now demonstrate these points.

6.2 Using Do-Calculus to Investigate the Roy Model

We show the limitations of the DoC for identifying the Roy model.



The first column of Table 20 presents the DAG of the original Roy model, which is denoted by G . The second column displays the DAG $G_{\underline{Z}}$ which suppresses the arrow arising from Z . The LMC of Z on DAG $G_{\underline{Z}}$ is $Z \perp\!\!\!\perp (Y, T)$. From Rule 2 in DoC, we obtain $P(T | do(Z)) = P(T | Z)$. Summarizing:

$$G_{\underline{Z}} \Rightarrow T \perp\!\!\!\perp Z \Rightarrow \text{by Rule 2 that } P(T | do(Z)) = P(T | Z). \quad (32)$$

Therefore, the modified Directed Acyclic Graph (DAG), $G_{\underline{Z}}$, enable us to assert that conditioning T on $Z = z$ is equivalent to choice T when we fix Z to value z . In the NR

framework, this result is obtained by the exogeneity condition $T(z) \perp\!\!\!\perp Z$, which states that the instrument Z is independent of the counterfactual choice $T(z)$ and thus $P(T|Z = z) = P(T(z))$ holds. Instrument Z in DAG $G_{\underline{Z}}$ is independent of both T and Y . This analysis also applies to Y . We can use (32) to obtain that $P(Y | do(Z)) = P(Y | Z)$, which means that conditioning on Z is equivalent to fixing Z . In summary, instrument Z is an external variable and the causal operation of fixing is translated to standard statistical conditioning.

The third column of Table 20 displays the DAG $G_{\overline{T}}$ which suppresses the arrow arriving at T . The LMC of Z in $G_{\overline{T}}$ implies $Z \perp\!\!\!\perp Y$. By Rule 1 of DoC, we have that $P(Y | do(T), Z) = P(Y | do(T))$. Summarizing:

$$G_{\overline{T}} \Rightarrow Y \perp\!\!\!\perp Z \Rightarrow \text{by Rule 1 that } P(Y | do(T), Z) = P(Y | do(T)). \quad (33)$$

This means that Z is statistically independent of Y when we fix T . This statement refers to the exogeneity condition $Y(t) \perp\!\!\!\perp Z$ or the independence relationship $Y \perp\!\!\!\perp Z | \tilde{T}$ of the HM framework.

The last column of Table 20 displays the DAG $G_{\overline{T}, \underline{Z}}$ which suppresses the arrow arriving at T and arising from Z . Note that the DAGs $G_{\overline{T}, \underline{Z}}$ and $G_{\overline{T}}$ are the same. The LMC of Z is $Z \perp\!\!\!\perp (T, Y, V)$ which implies that $Z \perp\!\!\!\perp T$ holds. Using Rule 2 of the DoC we obtain:

$$G_{\overline{T}, \underline{Z}} \Rightarrow Y \perp\!\!\!\perp Z | T \Rightarrow \text{by Rule 2 that } P(Y | do(T), do(Z)) = P(Y | do(T), Z). \quad (34)$$

Combining $P(Y | do(T), Z) = P(Y | do(T))$ in (33) with $P(Y | do(T), do(Z)) = P(Y | do(T), Z)$ in (34) we obtain $P(Y | do(T), do(Z)) = P(Y | do(T))$. This means that the probability distribution of the outcome Y when we fix both Z, T is the same as the counterfactual outcome generated by fixing only the choice T . In the NR framework, this property refers to the exclusion restriction $Y_i(t, z) = Y_i(t, z')$ for all $z, z' \in \text{supp}(Z)$.

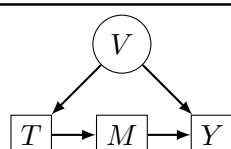
These statements *exhaust* the analysis of the Roy model analysis that can be performed using DoC. The method describes some key properties of the Roy model, but application of its rules alone cannot deliver identification of treatment effects. Indeed, the assumptions necessary for securing the identification of treatment effects in the Roy model cannot be

assessed by a DAG representation. Identifying assumptions, such as separability or monotonicity, impose restrictions on the functional form of the choice equation which go beyond the causal links described by a DAG.

6.3 The Front-door Model

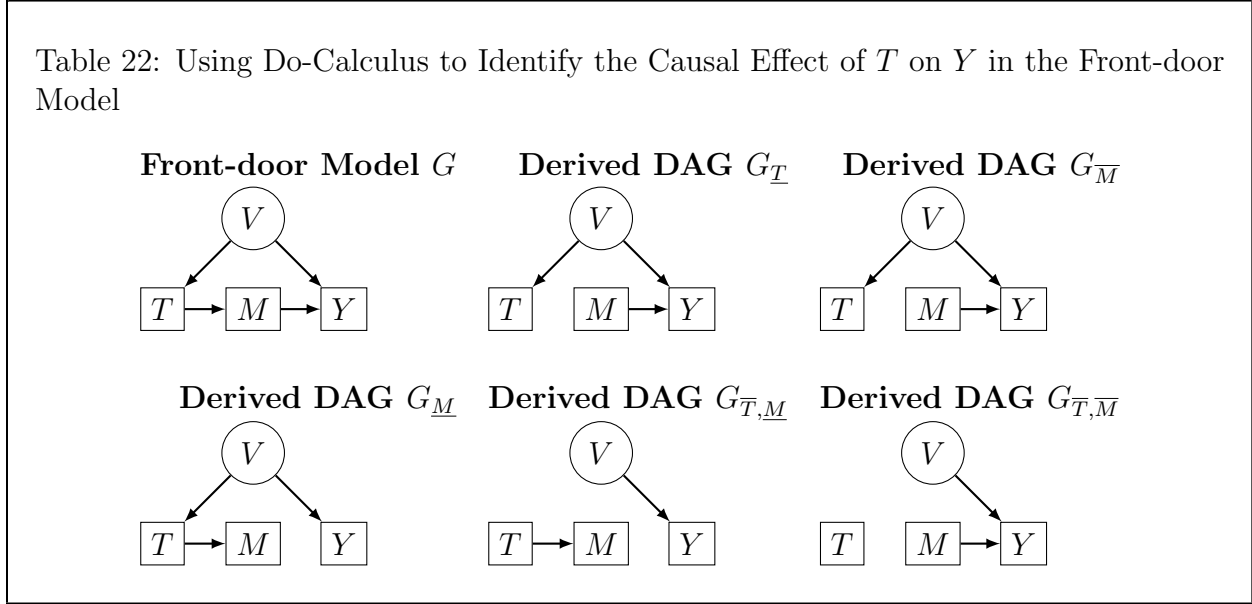
To make a more positive statement about do-calculus, it is useful to compare the identification machinery of the DoC and HM using a causal model when treatment effects *are* identified by DoC. We use the “Front-door model” of Pearl (2009b) to illustrate the differences in the approaches.

The Front-door model consists of three observed variables T, M, Y and an unobserved confounding variable V . Treatment T causes a mediator M which in turn causes outcome Y . Table 21 presents the causal representations of the model.

Table 21: Representations of the Front-door Model				
	Variable Map	Structural Eq.	DAG	LMC
V	$\mathbf{M}(V) = \emptyset$	$V = f_Z(\epsilon_V)$		$Z \perp\!\!\!\perp V \emptyset$
T	$\mathbf{M}(T) = \{V\}$	$T = f_T(V, \epsilon_T)$		$V \perp\!\!\!\perp \emptyset \emptyset$
M	$\mathbf{M}(M) = \{T\}$	$M = f_M(T, \epsilon_M)$		$M \perp\!\!\!\perp V T$
Y	$\mathbf{M}(Y) = \{M, V\}$	$Y = f_Y(M, V, \epsilon_Y)$		$Y \perp\!\!\!\perp \emptyset (M, V)$

The causal effect of T on Y in the Front-door model is identified. This result arises from the fact that the causal effect of T on M is not directly confounded by V since conditioning on T blocks the effect of the confounder V on M . Thus, we can identify the causal effect of M on Y conditional on T . The causal effect of T on Y can be evaluated as the compound effect of T on M and M on Y .

Table 22: Using Do-Calculus to Identify the Causal Effect of T on Y in the Front-door Model



We illustrate how to use DoC to identify the distribution of the counterfactual outcome $Y(t)$. To simplify notation, suppose that all variables are discrete. The do-calculus is cumbersome. The method requires the five derived DAGs displayed in Table 22. Identification for the counterfactual outcome is obtained by the following sequence of steps:

1. $T \perp\!\!\!\perp M$ in $G_{\underline{T}}$ holds, thus by Rule 2 we have that $P(M \mid do(T)) = P(M \mid T)$.
2. $M \perp\!\!\!\perp T$ in $G_{\overline{M}}$ holds, thus by Rule 3 we have that $P(T \mid do(M)) = P(T)$.
3. $M \perp\!\!\!\perp Y \mid T$ in $G_{\underline{M}}$ holds, thus by Rule 2 we have that $P(Y \mid T, do(M)) = P(Y \mid T, M)$.
4. Adding these results, we obtain:

$$\therefore P(Y \mid do(M)) = \sum_t P(Y \mid T = t, do(M))P(T = t \mid do(M))$$

by Law of Iterated Expectations (L.I.E.)

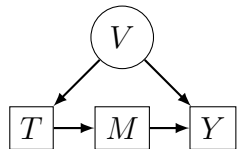
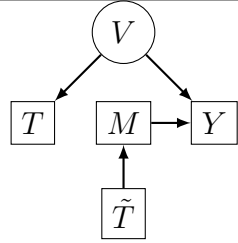
$$= \sum_t P(Y \mid T = t, M)P(T = t) \text{ by steps 1,2, and 3.}$$

5. $Y \perp\!\!\!\perp M \mid T$ in $G_{\overline{T}, \underline{M}}$ holds, thus by Rule 2, $P(Y \mid M, do(T)) = P(Y \mid do(M), do(T))$.
6. $Y \perp\!\!\!\perp T \mid M$ in $G_{\overline{T}, \overline{M}}$ holds, thus by Rule 3, $P(Y \mid do(T), do(M)) = P(Y \mid do(M))$.
7. Collecting these results, we have that $P(Y \mid Z, do(T)) = P(Y \mid do(Z), do(T)) = P(Y \mid do(M))$.
8. Finally, we can use previous results to obtain the following equation:

$$\begin{aligned}
\therefore P(Y \mid do(T) = t) &= \sum_m P(Y \mid M = m, do(T) = t)P(M = m \mid do(T) = t) \text{ by L.I.E.} \\
&= \sum_m P(Y \mid do(M) = m, do(T) = t)P(M = m \mid do(T) = t) \text{ by step 5.} \\
&= \sum_m P(Y \mid do(M) = m)P(M = m \mid do(T) = t) \text{ by step 7.} \\
&= \sum_m \left(\sum_{T=t'} P(Y \mid T=t', M=m)P(T=t') \right) P(M = m \mid T=t) \text{ by step 4.}
\end{aligned}$$

6.4 The Front Door Model in the Hypothetical Model Framework

We now investigate the same front-door model using the hypothetical framework. Table 24 displays the hypothetical model associated with the Front-door model as a DAG. The bottom panel of Table 24 presents the LMC for both models.

Table 23: The Empirical and Hypothetical Front-door Models	
Empirical Model	Hypothetical Model
	
LMC	LMC
$V \perp\!\!\!\perp M$	$V \perp\!\!\!\perp (M, \tilde{T})$
$T \perp\!\!\!\perp M \mid V$	$T \perp\!\!\!\perp (M, Y, \tilde{T}) \mid V$
$M \perp\!\!\!\perp V \mid T$	$M \perp\!\!\!\perp (T, V) \mid \tilde{T}$
$Y \perp\!\!\!\perp T \mid (V, M)$	$Y \perp\!\!\!\perp (T, \tilde{T}) \mid (V, M)$
	$\tilde{T} \perp\!\!\!\perp (T, V)$

Recall that counterfactual outcome in the hypothetical framework is denoted by the outcome Y conditioned on the hypothetical variable \tilde{T} . Identification consists of expressing the counterfactual outcome distribution $P_h(Y \mid \tilde{T} = t)$, which is defined in the hypothetical model, in terms of the observed distribution $P_e(T, M, Y)$, defined in the empirical model. The

connection between the probabilities of the hypothetical and empirical models is governed by the rules (16)–(17). The first rule states that, if $Y \perp\!\!\!\perp \tilde{T} \mid (T, W)$ holds for any variables Y, \tilde{T}, T, W in the hypothetical model, then we can equate $P_h(Y \mid \tilde{T} = t, T = t', W) = P_e(Y \mid T = t', W)$. On the other hand, $Y \perp\!\!\!\perp T \mid (\tilde{T}, W)$ implies that $P_h(Y \mid \tilde{T} = t, T = t', W) = P_h(Y \mid T = t, W)$.

The hypothetical framework requires analysts to find independence relationships of the hypothetical model that contain T and \tilde{T} . Useful relations are $Y \perp\!\!\!\perp \tilde{T} \mid (M, T)$ and $M \perp\!\!\!\perp T \mid \tilde{T}$.⁵⁴ It is also the case $T \perp\!\!\!\perp \tilde{T}$ holds as \tilde{T} is externally specified (exogenous) and does not cause T . We can then apply rules (16)–(17) to generate the following probability equalities:

$$Y \perp\!\!\!\perp \tilde{T} \mid (T, M) \quad \Rightarrow \quad P_h(Y \mid \tilde{T}, T = t', M) = P_e(Y \mid T = t', M). \quad (35)$$

$$M \perp\!\!\!\perp T \mid \tilde{T} \quad \Rightarrow \quad P_h(M \mid \tilde{T} = t, T) = P_e(M \mid T = t). \quad (36)$$

$$T \perp\!\!\!\perp \tilde{T} \mid T \quad \Rightarrow \quad P_h(T = t' \mid \tilde{T}) = P_e(T = t'). \quad (37)$$

The causal effect of T on Y of the Front-door model is identified through the following logic:

$$P_h(Y \mid \tilde{T} = t) = \sum_{t', m} P_h(Y \mid m, T = t', \tilde{T} = t) P_h(m \mid T = t', \tilde{T} = t) P_h(T = t' \mid \tilde{T} = t). \quad (38)$$

$$= \sum_{t', m} P_e(Y \mid m, T = t') P_e(m \mid T = t) P_e(T = t'). \quad (39)$$

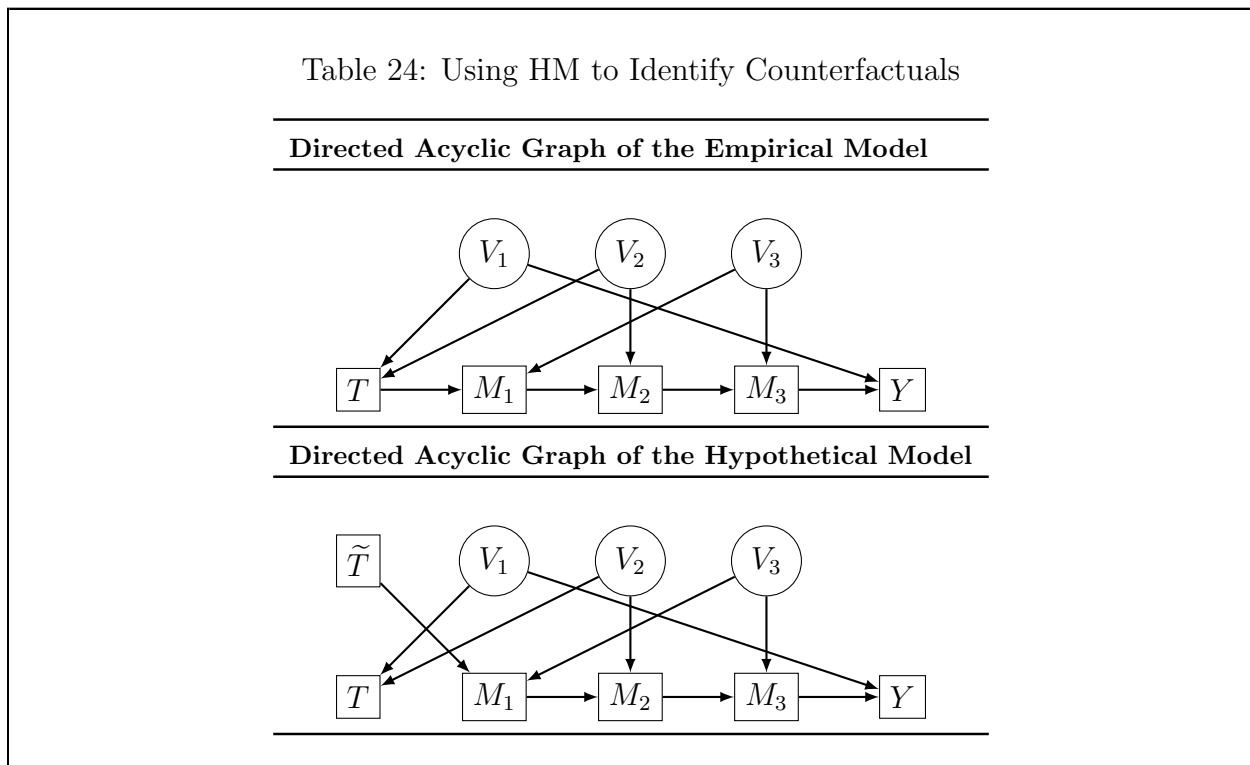
Equation (38) is a sum of probabilities defined in the hypothetical model by application of the law of iterated expectation over T and M . Equation (39) replaces each of the hypothetical model probabilities with empirical model probabilities listed in equations (35)–(37).

The identification of the counterfactual outcomes in the Front-door Model stems from the three independence relationships in (35)–(37). These independence relationships illustrate two properties that are at the core of the identification result. The first property is that the independence relationships alternate between T and \tilde{T} in the positions of conditioning

⁵⁴The first independence condition is due to the LMC $Y \perp\!\!\!\perp \tilde{T} \mid M$ and $(\tilde{T}, M) \perp\!\!\!\perp (T, V)$. The second one is due to the LMC of M .

variable and independent variable in the right-hand side. We term this property as *alternate conditionals*. The second property is that the sequence of conditioning variables on the right-hand side of (35)–(37) form a sequence $Y \rightarrow M \rightarrow T$ that starts at the targeted outcome Y and arrives at the treatment using the variable M to bridge these variables. We term this use of M as the *bridging* property.

Identification is secured whenever the properties of alternating conditionals and the bridging properties hold. We illustrate these ideas for the complex mediation model of Table 24. The model has three observed mediating variables M_1, M_2, M_3 (instead of M) and three unobserved, confounding variables V_1, V_2, V_3 (instead of V).



The following conditional independence relationships hold for the hypothetical model:

$$Y \perp\!\!\!\perp \tilde{T} \mid (T, M_3, M_2, M_1) \Rightarrow P_h(Y \mid \tilde{T}, T = t', M_3, M_2, M_1) = P_e(Y \mid T = t', M_3, M_2, M_1). \quad (40)$$

$$M_3 \perp\!\!\!\perp T \mid (\tilde{T}, M_2, M_1) \Rightarrow P_h(M_3 \mid \tilde{T} = t, T, M_2, M_1) = P_e(M_3 \mid T = t, M_2, M_1). \quad (41)$$

$$M_2 \perp\!\!\!\perp \tilde{T} \mid (T, M_1) \Rightarrow P_h(M_2 \mid \tilde{T}, T = t', M_1) = P_e(M_2 \mid T = t', M_1). \quad (42)$$

$$M_1 \perp\!\!\!\perp T \mid \tilde{T} \Rightarrow P_h(M_1 \mid \tilde{T} = t, T) = P_e(M_1 \mid T = t). \quad (43)$$

$$T \perp\!\!\!\perp \tilde{T} \mid T \Rightarrow P_h(T = t' \mid \tilde{T}) = P_e(T = t'). \quad (44)$$

The set of independence relationships (40)–(44) exhibits the alternate conditionals property. The first relationship is conditioned on T , the second on \tilde{T} , followed by T and so on. The bridging property also holds. The right-hand variable of each independence relationship provides a bridging sequence: $Y \rightsquigarrow M_3 \rightsquigarrow M_2 \rightsquigarrow M_1 \rightsquigarrow T$. The law of iterated expectations and independence relationships (40)–(44) enable us to express the counterfactual probability $P_h(Y \mid \tilde{T})$ as:

$$\textbf{Hypothetical Model} \quad P_h(Y \mid \tilde{T} = t) = \sum_{t', m_3, m_2, m_1} A_h \cdot B_h \cdot C_h \cdot D_h \cdot E_h.$$

where:

$$A_h = P_h(Y \mid m_3, m_2, m_1, T = t', \tilde{T} = t).$$

$$B_h = P_h(M_3 = m_3 \mid m_2, m_1, T = t', \tilde{T} = t).$$

$$C_h = P_h(M_2 = m_2 \mid m_1, T = t', \tilde{T} = t).$$

$$D_h = P_h(M_1 = m_1 \mid T = t', \tilde{T} = t).$$

$$E_h = P_h(T = t' \mid \tilde{T} = t).$$

Connection rules (16)–(17) enable us to translate hypothetical probabilities into the empirical probabilities as listed in (40)–(44). The resulting identification equation is presented below. It displays the alternative pattern of values t and t' in the same fashion as the identification equation of the Front-door model:

$$\textbf{Empirical Model} \quad P_e(Y(t)) = \sum_{t', m_3, m_2, m_1} A_e \cdot B_e \cdot C_e \cdot D_e \cdot E_e.$$

where:

$$A_e = P_e(Y \mid m_3, m_2, m_1, T = t').$$

$$B_e = P_e(M_3 = m_3 \mid m_2, m_1, T = t).$$

$$C_e = P_e(M_2 = m_2 \mid m_1, T = t').$$

$$D_e = P_e(M_1 = m_1 \mid T = t).$$

$$E_e = P_e(T = t').$$

6.5 Comparing DoC and HM Frameworks

Both DoC and HM employ structural equations and describe causal models with both observed and unobserved variables. They clearly separate the task of defining counterfactuals and identifying them. Both frameworks enable analysts to disentangle the tasks of causal analysis in Table 1. Both frameworks employ economic theory to define causal models (Task 1) and the structural equations that underlie the approach.

There are, however, some distinct practices in DoC and HM. When DoC fixes a treatment variable, it *eliminates* the equation for T in constructing the joint distribution of variables. All of the DoC analysis is done within the empirical model so generated.

HM does *not* eliminate the equation for the treatment variable. Instead, it *adds* a hypothetical variable. The presence of both treatment and hypothetical variables in the HM framework facilitates the study of the causal effects. HM can be used to readily analyze both external manipulation and conditioning, such as the treatment on the treated, whereas this is outside the scope of DoC. It facilitates examination of causal inference for direct and indirect effects in which the hypothetical variable replaces some but not all the treatment inputs of the structural equations. DoC invents new rules to undertake those tasks for each combination of conditioning variables.

The identification of causal effects (Task 2) requires connecting the hypothetical model with the empirical model. HM uses two statistical implications to connect the probability distributions of the hypothetical and empirical models. HM analyses remain within the realm of standard statistical theory and do not require invocation of non-probabilistic DAG-based rules.

The DoC machinery consists of three DAG-based rules. It constructs a series of possible DAGs. Each of them constitutes a causal model that modifies the empirical model. Each modification of the empirical model corresponds to introducing a new set of conditional independence relationships. The search for the combinations of DAGs and conditional inde-

pendence relationships that are required to identify counterfactuals grows exponentially. An algorithm has been developed to perform this task.⁵⁵ Calculations with HM are simpler than those based on DoC. They rely on a single modification of the original DAG, as encoded in the hypothetical model instead of a growing list of DAGs to implement the three guiding rules of DoC.

DoC relies critically on DAGs, conditional independence relationships, and a special set of rules. The HM machinery remains within the statistical realm to make statistics converse with causality. In doing so, the method is capable of accommodating assumptions that explore the rich variety of functional form restrictions, distributional assumptions, and cross-equation and cross-variable relationships that lie outside the scope of DoC.

7 Simultaneous Causality

The Generalized Roy model is usually expressed as a recursive model.⁵⁶ However, simultaneous causality is a property of many economic models. Examples of such models include those for social interactions, general equilibrium, Walrasian market clearing, or simultaneous play in models of industrial organization which are staples of economic theory (see, e.g., [Mas-Colell et al., 1995](#); [Tamer, 2003](#)). Such models are ignored in discussions of causality in the NR literature. The NR approach invokes the Stable Unit Treatment Value Assumption (SUTVA), which excludes the possibility of interactions among agents.⁵⁷ Such interactions are usually termed “confounders” and are treated as a problem rather than a source of information about economic and social behavior.

It is instructive to consider these models because they challenge the approaches used in the statistical literature, but are easily analyzed by rigorous econometric causal models. The pioneering econometric models studied by the Cowles Commission featured simultaneity.⁵⁸

⁵⁵See [Pearl \(2009b\)](#).

⁵⁶See, however, [Brock and Durlauf \(2007\)](#); [Heckman \(1978\)](#).

⁵⁷See, for instance, [Holland \(1986\)](#); [Imbens and Rubin \(2015\)](#).

⁵⁸See, e.g., [Koopmans et al. \(1950\)](#).

Haavelmo's (1943) paper explicitly analyzed causality in a simultaneous system. Many of the core ideas in simultaneous equations models are ignored or remain unknown to the followers of the statistical approaches, which rely on recursive formulations, which are considered to be essential features of causal models.

Simultaneous causality is an essential feature of many structural equation econometric models.⁵⁹ The LISREL model of Jöreskog (1973) allows for simultaneity, measurement error and latent variables proxied by measurements as discussed in Section 4.

The structural systems typically consist of two parts: (a) an autonomous structural system expressed in terms of latent variables (Bollen, 2002) and (b) a measurement system. The measurement system proxies the latent variables using observed measurements. The first part of the structural system consists of structure for person i :

$$\mathbf{B}\boldsymbol{\eta}_i = \boldsymbol{\alpha}_\eta + \boldsymbol{\Gamma}\boldsymbol{\chi}_i + \boldsymbol{\omega}_i \quad (45)$$

where $\boldsymbol{\omega}_i$, $\boldsymbol{\eta}_i$, $\boldsymbol{\chi}_i$ are vectors of latent variables. The measurement system consists of vectors of measurements:

$$\text{Measurement: } \begin{cases} \mathbf{y}_i = \boldsymbol{\alpha}_y + \boldsymbol{\Lambda}_y\boldsymbol{\eta}_i + \boldsymbol{\varepsilon}_i & (\text{measurement for } \eta_i) \\ \mathbf{x}_i = \boldsymbol{\alpha}_x + \boldsymbol{\Upsilon}_x\boldsymbol{\chi}_i + \boldsymbol{\xi}_i & (\text{measurement for } \chi_i) \end{cases}$$

where $\boldsymbol{\varepsilon}_i$, $\boldsymbol{\Upsilon}_i$, and $\boldsymbol{\xi}_i$ are vectors of latent variables. These models have been extended to time series and panel data settings (see e.g. Bollen, 1989; Goldberger and Duncan, 1973; Hansen and Sargent, 1982).

In a valuable paper, Bollen and Pearl (2013) exposit this system of equations as a causal model with simultaneity and show how various measurement systems use factor models and other approaches to proxy the latent variables which may be the variables measured with error or omitted variables, like ability in an earnings equation, or technical efficiency in

⁵⁹See Goldberger (1972); Haavelmo (1943, 1944); Koopmans et al. (1950), Goldberger and Duncan (1973).

a production function. They dispel many misguided criticisms of the structural approach lodged by advocates of the NR approach. These systems are equipped to use cross equation restrictions and covariance restrictions to secure identification of causal parameters. [Hansen and Sargent \(1982\)](#) is an example of this approach applied to time series models.

This literature is rich and we lack the space to exposit it thoroughly. We note that as previously discussed linear equation versions of these models provide a framework for proxying V . It is also an approach for studying mediation where analysts can study how interventions on χ_i percolate through equation system (43). [Schennach \(2020\)](#) summarizes a large literature on nonparametric factor and proxy models.

Instead of a general exposition of these systems, we refer the reader to [Bollen and Pearl \(2013\)](#) and consider a simple two-equation simultaneous equations model of the sort exposted by [Haavelmo \(1943\)](#). We consider a system of two autonomous (structurally-invariant) causal equations:

$$Y_1 = g_{Y_1}(Y_2, X_1, U_1, \epsilon_1) \tag{46}$$

$$Y_2 = g_{Y_2}(Y_1, X_2, U_2, \epsilon_2) \quad U_1 \not\perp U_2. \tag{47}$$

We use this system to demonstrate how causality can be analyzed in simultaneous systems. Again, ϵ_1 and ϵ_2 are mutually independent and independent of U_1, U_2, X_1 , and X_2 .

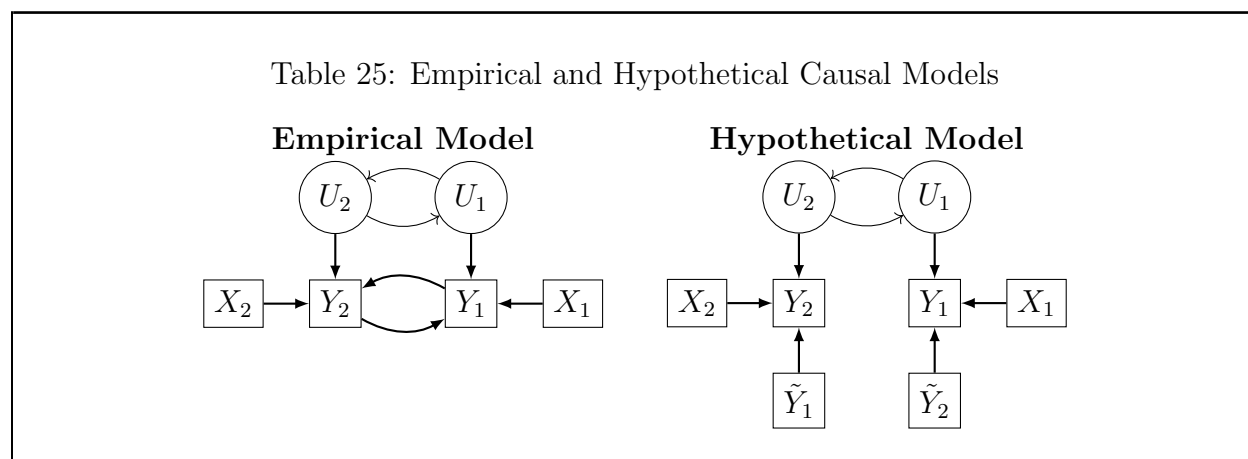
This system of equations represents two maps: $g_{Y_1}: (Y_2, X_1, U_1) \rightarrow Y_1; g_{Y_2}: (Y_1, X_2, U_2) \rightarrow Y_2$. Y_1 and Y_2 could be actions of a pair of interacting agents.⁶⁰ To simplify the discussion, we assume that both equations are twice continuously differentiable. This is a convenience and not a necessity. The model of equations (46)–(47) is treated in a special way in the DoC approach. Models with multiple simultaneous equations are standard in the literature (see, e.g., [Bollen, 1989](#); [Fisher, 1966b](#); [Goldberger and Duncan, 1973](#); [Koopmans et al., 1950](#); [Theil, 1958, 1971](#)).

⁶⁰In the literature on peer effects, simultaneous equation problems are relabeled as “reflection problems.” See [Manski \(1993\)](#); [Moffitt \(2001\)](#).

As previously noted, equations (46) and (47) are assumed to be structural, i.e., invariant under manipulations of their arguments, so they are stable, autonomous maps. Policies consist of manipulations of their arguments. Autonomy is one part of the *SUTVA* assumption in the NR model.⁶¹

In the classical model of market clearing equilibrium, Y_1 is price; Y_2 is quantity and X_1 , X_2 , U_1 , and U_2 are causal determinants. Equations (46) and (47) are generated by thought experiments varying the arguments and tracing out the outcomes. Thus, Y_1 in (46) could be the market price that is consistent with hypothetical values Y_2, X_1, U_1 . (47) is the analogous relationship for market quantity Y_2 . The addition of unobserved (by the economist) variables U_1 and U_2 is made in anticipation of empirical applications. In the peer effects literature, Y_1 and Y_2 are behaviors of two interacting agents (e.g., smoking or drug use).

In terms of our previous notation, the variable set for the empirical model is $\mathcal{T}_e = \{Y_1, Y_2, X_1, X_2, U_1, U_2\}$. $\mathbb{M}_e(Y_1) = \{Y_2, X_1, U_1\}$ and $\mathbb{M}_e(Y_2) = \{Y_1, X_2, U_2\}$. The empirical and hypothetical models are displayed as DAGs in Table 25 given by:⁶²



The LMC condition does not apply so that the Bayesian net approach fails. “Fixing” and the hypothetical model approach readily extend to a system of simultaneous equations for Y_1 and Y_2 , whereas the fundamentally recursive methods based on DAGs require special treatment.

⁶¹The other part is absence of simultaneity or general equilibrium effects. See Heckman (2008a).

⁶² U_2 and U_1 are reciprocally related.

7.1 Completeness

“Completeness” assumes the existence of at least a local solution for Y_1 and Y_2 in terms of (X_1, X_2, U_1, U_2) :

$$Y_1 = \phi_1(X_1, X_2, U_1, U_2) \quad (48)$$

$$Y_2 = \phi_2(X_1, X_2, U_1, U_2). \quad (49)$$

These are *reduced form* equations (see, e.g., [Koopmans et al., 1950](#); [Matzkin, 2008, 2013](#)). They inherit the autonomy properties of the structural equations. Completeness is a property that guarantees the conceptual possibility of simultaneity, which is not necessarily guaranteed. If it fails, the existence of consistent solutions to (46) and (47) is not guaranteed. Nonetheless autonomous correspondences may still exist and they can be used to make set-valued causal inferences.⁶³

The causal effect of Y_2 on Y_1 when Y_2 is fixed at y_2 is generated by

$$Y_1(y_2) = g_{Y_1}(y_2, X, U_1).$$

Symmetrically, the causal effect of Y_1 on Y_2 when Y_1 is fixed at y_1 is generated by:

$$Y_2(y_1) = g_{Y_2}(y_1, X, U_2).$$

The relationships (46) and (47) can be defined even if they might not be identified or estimated. The *completeness assumption* says that there are values of X_1, X_2, U_1, U_2 that generate values of Y_1, Y_2 consistent with (46) and (47). These involve hypothetical variations. For certain models no such sets of variables may exist and the models are termed incomplete.

7.2 Can We Hypothetically Vary Y_2 and Y_1 ?

If Y_2 and Y_1 are simultaneously determined, the notion of varying Y_2 to change Y_1 may seem impossible. [Pearl \(2009a\)](#) maintains his focus on recursive models and addresses this problem in a very special way by assuming structural invariance and “shutting one equation down,”

⁶³See, e.g., [Heckman \(1978\)](#); [Mas-Colell et al. \(1995\)](#); [Quandt \(1988\)](#); [Tamer \(2003\)](#).

assuming the rest of the system remains unchanged. Thus, for example, equation (47) is suspended, but (46) is maintained. This is consistent with the logic of do-calculus, which eliminates relationships from systems, assuming invariance of the remaining system. He sets Y_2 to a constant that can be manipulated in (46). This thought experiment converts a simultaneous system into a recursive system with all other equations assumed to hold as before. After Y_2 is fixed, the do-calculus can be applied.

This approach is cumbersome and strains credibility in many economic contents (e.g., person 1 influences 2, but not vice versa), but it is logically possible.⁶⁴ It is unnecessary if exclusions in (46) and (47) are used. To show this, we define exclusion of X_2 in (46) as $\frac{\partial g_{Y_1}}{\partial X_2} = 0$ for all (Y_2, X_1, X_2, U_1) .⁶⁵ Exclusion of X_1 in (47) is defined as $\frac{\partial g_{Y_2}}{\partial X_1} = 0$ for all (Y_1, X_1, X_2, U_2) . Implicit is the assumption that components of X_1 and X_2 can be varied. Under completeness and exclusion X_2 from (47), by the chain rule, the causal effect of Y_2 on Y_1 is

$$\frac{\partial g_{Y_1}}{\partial Y_2} = \frac{\partial Y_1}{\partial X_2} \Big/ \frac{\partial Y_2}{\partial X_2} = \frac{\partial \theta_1}{\partial X_2} \Big/ \frac{\partial \theta_2}{\partial X_2}.$$

We may define and identify the causal effects for Y_1 on Y_2 in an analogous fashion. Variations in X_1 and X_2 that respect completeness define the causal parameters when the components of X_1 and X_2 can be independently varied.⁶⁶ No implausible “shutting down” of any equation in a system while assuming autonomy (structural invariance) of the remaining system is required.

This logic is now standard and is the basis for an estimation technique, “indirect least squares” (see [Theil, 1958](#) and [Tinbergen, 1930, 1939](#)). It demonstrates the flexibility of the econometric approach for defining and identifying causal parameters outside the narrow world of DAGs. [Fisher \(1966b\)](#) gives a range of approaches for identifying systems like

⁶⁴In a market for a good, shutting down the supply equation would likely alter the properties of demand curves as agents would face a different market structure altering their expectations. Construction of a theory consistent counterfactual world would entail such considerations.

⁶⁵Or more generally, X_2 is not an argument of g_{Y_1} .

⁶⁶The completeness condition is part of the hypothetical model thought experiment. In some contexts it may be ruled out as not credible.

(46) and (47) and more general versions using restrictions within and across equations for observables and unobservables.

7.3 Econometric Mediation Analysis

We have already discussed mediation analyses in recursive models. These notions extend to models with simultaneity. Under completeness, reduced form (48) characterizes the **net effect** of a policy change X_1 :

$$\frac{\partial Y_1}{\partial X_1} = \frac{\partial \phi_1(X_1, X_2, U_1, U_2)}{\partial X_1}. \quad (50)$$

Following Klein and Goldberger (1955) and Wright (1921, 1934), we can conduct “mediation analyses” that address problem **P-2** and trace the impact of an externally manipulated X_1 on Y_1 , both through its direct effect on (46) and its indirect effect through Y_2 :

$$\frac{\partial Y_1}{\partial X_1} = \underbrace{\left(\frac{\partial g_{Y_1}}{\partial Y_2} \right)}_{\text{From Structure}} \underbrace{\left(\frac{\partial Y_2}{\partial X_1} \right)}_{\text{From Reduced Form}} + \underbrace{\frac{\partial g_{Y_1}}{\partial X_1}}_{\text{From Structure}} = \frac{\partial \phi_1(X_1, X_2, U_1, U_2)}{\partial X_1}$$

Indirect effect through Y_2 Direct effect

This approach can be readily applied to recursive systems and general multiple equation systems. Reliance on linear equations, while traditional in the literature, is not necessary and nonparametric approaches are available.⁶⁷

Mediation analysis is a staple of econometric policy evaluation to examine all channels of influence of variables (see, e.g., Theil, 1958). All of the tools used to analyze simultaneous equations are available to estimate these models (See e.g., Amemiya, 1985; Fisher, 1966b; Matzkin, 2007). Klein and Goldberger (1955) is a classic example of dynamic mediation analysis in a Keynesian model of the time series of consumption and investment in the U.S. economy.

⁶⁷See Matzkin (2008, 2013, 2015) for nonparametric analyses of such systems.

8 Conclusion

This paper presents the basic framework of the econometric model for causal analysis. We discuss the definition of causal parameters and approaches to their identification within it. We consider two recent approaches to causality that are used in the non-economic literature on causal inference and their relationship with the econometric approach.

The econometric model is based on clearly stated and interpretable models of behavior that characterize the lessons of economic theory and allow for testing it, for synthesizing evidence on it from multiple studies, constructing credible policy counterfactuals, including forecasting policy impacts in new environments and forecasting the likely impacts of policies never previously implemented. The econometric approach delineates the definition of causal parameters, their identification and their estimation as three separate tasks.

The two competing statistical approaches are: (a) the Neyman-Rubin (NR) approach rooted in the statistics of experiments, and (b) the *do-calculus* (DoC) that originated in computer science. Both address some of the same problems tackled by the econometric approach. Each has important, but different, limitations. Neither has the flexibility or clarity of the econometric approach.

All start from the basic intuitive definition of a causal effect as a *ceteris paribus* consequence of a change in inputs on outcomes, where the change can be a policy. However, the rules for constructing and identifying counterfactuals are very different in these approaches.

The do-calculus (DoC) invokes a special set of rules for identifying causal parameters that lie outside of probability theory and that use a limited class of identifying assumptions for behavioral equations. It relies heavily on recursive directed-acyclic-graphs and assumptions about conditional independence relationships. Its rigid rules preclude the use of many traditional techniques of identification and estimation.

The Neyman-Rubin (NR) approach eschews the benefits of structural equations and many fruitful strategies for their identification. Reflecting its origins, it casts all policy problems into a “treatment-control” framework. Randomized experiments rather than thought experiments are foundational elements in this approach. In some versions, it conflates issues of definition with issues of identification. Its lack of reliance on structural equations with explicit links to theory and explicit analyses of unobservables, makes it difficult to interpret estimates obtained from it to analyze well-posed economic questions with it using the large toolkit of modern econometrics, or to synthesize studies within a common framework.

Econometrics has a rich body of theory and tools to address policy problems. Applied economists would do well by using the impressive set of conceptual tools available from econometric theory.

References

- Aakvik, A., J. J. Heckman, and E. J. Vytlacil (1999). Training effects on employment when the training effects are heterogeneous: An application to Norwegian vocational rehabilitation programs. University of Bergen Working Paper 0599, and University of Chicago.
- 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(1–2), 15–51.
- Abadie, A. (2005). Semiparametric difference-in-differences estimators. *The Review of Economic Studies* 72(1), 1–19.
- 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. In J. J. Heckman and E. E. Leamer (Eds.), *Handbook of Econometrics*, Volume 6B, Chapter 72, pp. 5145–5303. Amsterdam: Elsevier Science B. V.
- Altonji, J. G. and R. L. Matzkin (2005, July). Cross section and panel data estimators for nonseparable models with endogenous regressors. *Econometrica* 73(4), 1053–1102.
- Amemiya, T. (1985). *Advanced Econometrics*. Cambridge, MA: Harvard University Press.
- Angrist, J. D., G. W. Imbens, and D. Rubin (1996). Identification of causal effects using instrumental variables. *Journal of the American Statistical Association* 91(434), 444–455.

- Angrist, J. D. and J.-S. Pischke (2009). *Mostly Harmless Econometrics: An Empiricist's Companion*. Princeton, NJ: Princeton University Press.
- Bareinboim, E. and J. Pearl (2016). Causal inference and the data-fusion problem. *Proceedings of the National Academy of Sciences* 113(27), 7345–7352.
- Bertrand, M., E. Duflo, and S. Mullainathan (2004, February). How much should we trust differences-in-differences estimates? *Quarterly Journal of Economics* 119(1), 249–275.
- Bjerkholt, O. and A. Dupont (2010). Ragnar Frisch's conception of econometrics. *History of Political Economy* 42(1), 21–73.
- Blundell, R., A. Duncan, and C. Meghir (1998, July). Estimating labor supply responses using tax reforms. *Econometrica* 66(4), 827–861.
- Blundell, R. and J. Powell (2003). Endogeneity in nonparametric and semiparametric regression models. In L. P. H. M. Dewatripont and S. J. Turnovsky (Eds.), *Advances in Economics and Econometrics: Theory and Applications, Eighth World Congress*, Volume 2. Cambridge, UK: Cambridge University Press.
- Bollen, K. A. (1989). *Structural Equations with Latent Variables*. New York: Wiley.
- Bollen, K. A. (2002). Latent variables in psychology and the social sciences. *Annual Review of Psych* 53(1), 605–634.
- Bollen, K. A. and J. Pearl (2013). Eight myths about causality and structural equation models. In S. L. Morgan (Ed.), *Handbook of Causal Analysis for Social Research*, Chapter 15, pp. 301–328. Springer, Dordrecht.
- Brock, W. A. and S. N. Durlauf (2007). Identification of binary choice models with social interactions. *Journal of Econometrics* 140(1), 52–75.
- Buchinsky, M. and R. Pinto (2021). Using economic incentives to generate monotonicity criteria of iv models. Unpublished Manuscript, UCLA.
- Bursztyn, L. and D. Y. Yang (2021). Misperceptions about others. Working Paper 29168, NBER. Unpublished.
- Cameron, S. V. and J. J. Heckman (1998, April). Life cycle schooling and dynamic selection bias: Models and evidence for five cohorts of American males. *Journal of Political Economy* 106(2), 262–333.
- Carneiro, P., K. Hansen, and J. J. Heckman (2003, May). Estimating distributions of treatment effects with an application to the returns to schooling and measurement of the effects of uncertainty on college choice. *International Economic Review* 44(2), 361–422.
- Chatfield, C. (2000). *Time-Series Forecasting*. CRC Press.
- Cox, D. R. (1958). *Planning of Experiments*. New York: Wiley.

- Cunha, F., J. J. Heckman, and S. Navarro (2005, April). Separating uncertainty from heterogeneity in life cycle earnings, The 2004 Hicks Lecture. *Oxford Economic Papers* 57(2), 191–261.
- Cunha, F., J. J. Heckman, and S. M. Schennach (2010, May). Estimating the technology of cognitive and noncognitive skill formation. *Econometrica* 78(3), 883–931.
- Dawid, A. P. (1976). Properties of diagnostic data distributions. *Biometrics* 32(3), 647–658.
- Dawid, A. P. (1979). Conditional independence in statistical theory (with discussion). *Journal of the Royal Statistical Society. Series B (Statistical Methodological)* 41(1), 1–31.
- Dippel, C., R. Gold, S. Heblich, and R. Pinto (2020). Mediation analysis in iv settings with a single instrument. *Unpublished Manuscript*.
- Ekeland, I., J. J. Heckman, and L. Nesheim (2004, February). Identification and estimation of hedonic models. *Journal of Political Economy* 112(S1), S60–S109. Paper in Honor of Sherwin Rosen: A Supplement to Volume 112.
- Fisher, F. (1966a). *The Identification Problem in Econometrics*. Economics handbook series. McGraw-Hill.
- Fisher, F. M. (1966b). *The Identification Problem in Econometrics*. New York: McGraw-Hill.
- Fisher, R. A. (1935). *The Design of Experiments*. London: Oliver and Boyd.
- Frangakis, C. E. and D. Rubin (2002). Principal stratification in causal inference. *Biometrics* 58(1), 21–29.
- Frisch, R. (1930). A dynamic approach to economic theory: Lectures by Ragnar Frisch at Yale University. Lectures at Yale University beginning September, 1930. Mimeographed, 246 pp. Frisch Archives, Department of Economics, University of Oslo.
- Frisch, R. (1933). Problèmes et méthodes de l'économétrie. Eight lectures given at Institut Henri Poincaré, University of Paris, March–April 1933. Frisch Archive, Department of Economics, University of Oslo.
- Frisch, R. (1933, published 2009). In O. Bjerkholt and A. Dupont-Kieffer (Eds.), *Problems and Methods of Econometrics: The Poincaré Lectures of Ragnar Frisch, 1933*. New York, New York: Routledge.
- Frisch, R. (1938). Autonomy of economic relations: Statistical versus theoretical relations in economic macrodynamics. Paper given at League of Nations. Reprinted in D.F. Hendry and M.S. Morgan (1995), *The Foundations of Econometric Analysis*, Cambridge University Press.
- Geiger, D., T. Verma, and J. Pearl (1990). Identifying independence in bayesian networks. *Networks* 20(5), 507–534.

- Glymour, C., R. Scheines, and P. Spirtes (2014). *Discovering Causal Structure: Artificial Intelligence, Philosophy of Science, and Statistical Modeling*. Academic Press.
- Goldberger, A. S. (1972, November). Structural equation methods in the social sciences. *Econometrica* 40(6), 979–1001.
- Goldberger, A. S. and O. D. Duncan (1973). Structural equation models in the social sciences. In O. D. Duncan and A. S. Goldberger (Eds.), *Social Science Research Council (États-Unis) and University of Wisconsin. Social Systems Research Institute*. New York: Seminar Press.
- Greenland, S., J. Pearl, and J. Robins (1999). Causal diagrams for epidemiologic research. *Epidemiology* 10 1, 37–48.
- Haavelmo, T. (1943, January). The statistical implications of a system of simultaneous equations. *Econometrica* 11(1), 1–12.
- Haavelmo, T. (1944). The probability approach in econometrics. *Econometrica* 12(Supplement), iii–vi and 1–115.
- Hamilton, J. D. (2000). *Time Series Analysis*. Princeton University Press.
- Hansen, B. E. (2022). *Econometrics*. Princeton University Press. Princeton University Press.
- Hansen, L. P. and T. J. Sargent (1982, May). Instrumental variables procedures for estimating linear rational expectations models. *Journal of Monetary Economics* 9(3), 263–296.
- Heckman, J. and R. Pinto (2018). Unordered monotonicity. *Econometrica* 86, 1–35.
- Heckman, J. J. (1978, July). Dummy endogenous variables in a simultaneous equation system. *Econometrica* 46(4), 931–959.
- Heckman, J. J. (1979, January). Sample selection bias as a specification error. *Econometrica* 47(1), 153–162.
- Heckman, J. J. (2008a, April). Econometric causality. *International Statistical Review* 76(1), 1–27.
- Heckman, J. J. (2008b). The principles underlying evaluation estimators with an application to matching. *Annales d’Economie et de Statistiques* 91–92, 9–73.
- Heckman, J. J., H. Ichimura, J. Smith, and P. E. Todd (1998, September). Characterizing selection bias using experimental data. *Econometrica* 66(5), 1017–1098.
- Heckman, J. J., R. J. LaLonde, and J. A. Smith (1999). The economics and econometrics of active labor market programs. In O. C. Ashenfelter and D. Card (Eds.), *Handbook of Labor Economics*, Volume 3A, Chapter 31, pp. 1865–2097. New York: North-Holland.
- Heckman, J. J. and E. E. Leamer (2001). *Handbook of Econometrics*, Volume 5 of *Handbooks in Economics*. Amsterdam: North Holland.

- Heckman, J. J. and E. E. Leamer (2007). *Handbook of Econometrics*, Volume 6AB of *Handbooks in Economics*. Amsterdam: North Holland.
- Heckman, J. J. and S. Navarro (2004, February). Using matching, instrumental variables, and control functions to estimate economic choice models. *Review of Economics and Statistics* 86(1), 30–57.
- Heckman, J. J. and R. Pinto (2015). Causal analysis after Haavelmo. *Econometric Theory* 31(1), 115–151.
- Heckman, J. J. and R. Robb (1985a). Alternative methods for evaluating the impact of interventions. In J. J. Heckman and B. S. Singer (Eds.), *Longitudinal Analysis of Labor Market Data*, Volume 10, pp. 156–245. New York: Cambridge University Press.
- Heckman, J. J. and R. Robb (1985b, October–November). Alternative methods for evaluating the impact of interventions: An overview. *Journal of Econometrics* 30(1–2), 239–267.
- Heckman, J. J. and B. S. Singer (1984, March). A method for minimizing the impact of distributional assumptions in econometric models for duration data. *Econometrica* 52(2), 271–320.
- Heckman, J. J. and C. Taber (2008). The roy model. In S. N. Durlauf and L. E. Blume (Eds.), *New Palgrave Dictionary of Economics* (2 ed.). Basingstoke, UK: Palgrave Macmillan.
- Heckman, J. J., S. Urzúa, and E. Vytlavil (2008). Instrumental variables in models with multiple outcomes: the general unordered case. *Annales d'économie et de Statistique* (91/92), 151–174.
- Heckman, J. J. and E. J. Vytlacil (1999, April). Local instrumental variables and latent variable models for identifying and bounding treatment effects. *Proceedings of the National Academy of Sciences* 96(8), 4730–4734.
- Heckman, J. J. and E. J. Vytlacil (2005, May). Structural equations, treatment effects and econometric policy evaluation. *Econometrica* 73(3), 669–738.
- Heckman, J. J. and E. J. Vytlacil (2007a). Econometric evaluation of social programs, part I: Causal models, structural models and econometric policy evaluation. In J. J. Heckman and E. E. Leamer (Eds.), *Handbook of Econometrics*, Volume 6B, Chapter 70, pp. 4779–4874. Amsterdam: Elsevier B. V.
- Heckman, J. J. and E. J. Vytlacil (2007b). Econometric evaluation of social programs, part II: Using the marginal treatment effect to organize alternative economic estimators to evaluate social programs, and to forecast their effects in new environments. In J. J. Heckman and E. E. Leamer (Eds.), *Handbook of Econometrics*, Volume 6B, Chapter 71, pp. 4875–5143. Amsterdam: Elsevier B. V.
- Holland, P. W. (1986, December). Statistics and causal inference. *Journal of the American Statistical Association* 81(396), 945–960.

- Holland, P. W. (1997). Some reflections on Freedmans critiques. In *Topics in the Foundation of Statistics*, pp. 50–57. Springer.
- Hoyer, P. O., D. Janzing, J. M. Mooij, J. Peters, and B. Schölkopf (2009). Nonlinear causal discovery with additive noise models. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou (Eds.), *Advances in Neural Information Processing Systems 21*, pp. 689–696. Curran Associates, Inc.
- Huang, Y. and M. Valtorta (2006). Pearl’s calculus of intervention is complete. In *Proceedings of the Twenty-Second Conference on Uncertainty in Artificial Intelligence, UAI’06*, Arlington, Virginia, USA, pp. 217224. AUAI Press.
- Hurwicz, L. (1962). On the structural form of interdependent systems. In E. Nagel, P. Suppes, and A. Tarski (Eds.), *Logic, Methodology and Philosophy of Science*, pp. 232–239. Stanford University Press.
- Imai, K., L. Keele, D. Tingley, and T. Yamamoto (2011). Unpacking the black box of causality: Learning about causal mechanisms from experimental and observational studies. *American Political Science Review* 105, 765–789.
- Imai, K., L. Keele, and T. Yamamoto (2010). Identification, inference and sensitivity analysis for causal mediation effects. *Statistical Science* 25(1), 51–71.
- Imbens, G. W. and J. D. Angrist (1994, March). Identification and estimation of local average treatment effects. *Econometrica* 62(2), 467–475.
- Imbens, G. W. and D. B. Rubin (2015). *Causal Inference for Statistics, Social, and Biomedical Sciences: An Introduction*. Cambridge University Press.
- Jöreskog, K. G. (1973). Analysis of covariance structures. In *Multivariate Analysis–III*, pp. 263–285. Elsevier.
- Kiiveri, H., T. P. Speed, and J. B. Carlin (1984). Recursive causal models. *Journal of the Australian Mathematical Society (Series A)*, 30–52.
- Klein, L. R. and A. S. Goldberger (1955). *An Econometric Model of the United States, 1929–1952*. Amsterdam: North-Holland Publishing Company.
- Knight, F. (1921). *Risk, Uncertainty and Profit*. New York: Houghton Mifflin Company.
- Koopmans, T. C., H. Rubin, and R. B. Leipnik (1950). Measuring the equation systems of dynamic economics. In T. C. Koopmans (Ed.), *Statistical Inference in Dynamic Economic Models*, Number 10 in Cowles Commission Monograph, Chapter 2, pp. 53–237. New York: John Wiley & Sons.
- Lauritzen, S. L. (1996). *Graphical Models*. Oxford, UK: Clarendon Press.
- Lee, S. and B. Salanié (2018). Identifying effects of multivalued treatments. *Econometrica* 86, 1939–1963.

- Levin, B. and H. Robbins (1983, Autumn). Urn models for regression analysis, with applications to employment discrimination studies. *Law and Contemporary Problems* 46(4, Statistical Inference in Litigation), 247–267.
- Lewbel, A. (2019). The identification zoo: Meanings of identification in econometrics. *Journal of Economic Literature* 57(4), 835–903.
- Lopez-Paz, D., R. Nishihara, S. Chintala, B. Schölkopf, and L. Bottou (2017). Discovering causal signals in images. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 6979–6989.
- Lucas, Jr., R. E. (1976). Econometric policy evaluation: A critique. In K. Brunner and A. H. Meltzer (Eds.), *The Phillips Curve and Labor Markets*, Volume 1 of *Carnegie-Rochester Conference Series on Public Policy*. Amsterdam: North-Holland.
- Manski, C. F. (1993, July). Identification of endogenous social effects: The reflection problem. *Review of Economic Studies* 60(3), 531–542.
- Marschak, J. (1953). Economic measurements for policy and prediction. In W. C. Hood and T. C. Koopmans (Eds.), *Studies in Econometric Method*, pp. 1–26. New Haven, CT: Yale University Press.
- Marshall, A. (1961). *Principles of Economics* (Ninth (Valorium) Edition ed.). London, Macmillan for the Royal Economic Society.
- Mas-Colell, A., M. D. Whinston, and J. R. Green (1995). *Microeconomic Theory*. New York: Oxford University Press.
- Matzkin, R. L. (1993, July). Nonparametric identification and estimation of polychotomous choice models. *Journal of Econometrics* 58(1–2), 137–168.
- Matzkin, R. L. (2007). Nonparametric identification. In J. J. Heckman and E. E. Leamer (Eds.), *Handbook of Econometrics*, Volume 6B. Amsterdam: Elsevier.
- Matzkin, R. L. (2008). Identification in nonparametric simultaneous equations models. *Econometrica* 76(5), 945–978.
- Matzkin, R. L. (2013). Nonparametric identification of structural economic models. *Annual Review of Economics* 5(1), 457–486.
- Matzkin, R. L. (2015). Estimation of nonparametric models with simultaneity. *Econometrica* 83(1), 1–66.
- McFadden, D. (1974). Conditional logit analysis of qualitative choice behavior. In P. Zarembka (Ed.), *Frontiers in Econometrics*, pp. 105–142. New York: Academic Press.
- Moffitt, R. A. (2001). Policy interventions, low-level equilibria, and social interactions. In S. Durlauf and P. Young (Eds.), *Social Dynamics*, Volume 4, pp. 6–17. MIT Press.

- Mogstad, M. and A. Torgovitsky (2018). Identification and extrapolation of causal effects with instrumental variables. *Annual Review of Economics* 2, 577–613.
- Morgan, S. L. and C. Winship (2015). *Counterfactuals and Causal Inference*. Cambridge University Press.
- Nerlove, M. (1967). Recent empirical studies of the CES and related production functions. In *The Theory and Empirical Analysis of Production*, pp. 55–136. National Bureau of Economic Research.
- Neyman, J. (1923). Statistical problems in agricultural experiments. *Journal of the Royal Statistical Society II (Supplement)*(2), 107–180.
- Olley, G. S. and A. Pakes (1996, November). The dynamics of productivity in the telecommunications equipment industry. *Econometrica* 64(6), 1263–1297.
- Pearl, J. (1988). *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc.
- Pearl, J. (1993). [Bayesian analysis in expert systems]: Comment: Graphical models, causality and intervention. *Statistical Science* 8(3), 266–269.
- Pearl, J. (1995, December). Causal diagrams for empirical research. *Biometrika* 82(4), 669–688.
- Pearl, J. (2009a). Causal inference in statistics: An overview. *Statistics Surveys* 3, 96–146.
- Pearl, J. (2009b). *Causality: Models, Reasoning, and Inference* (2nd ed.). New York: Cambridge University Press.
- Pearl, J. (2009c). Myth, confusion, and science in causal analysis. *Technical Report, UCLA, Department of Statistics*.
- Pearl, J. (2012). The do-calculus revisited. *CoRR abs/1210.4852*.
- Peters, J., D. Jazzing, and B. Schölkopf (2017). *Elements of Causal Inference: Foundations and Learning Algorithms*. Cambridge, MA: MIT Press.
- Prakasa Rao, B. L. S. (1992). *Identifiability in Stochastic Models: Characterization of Probability Distributions*. Probability and mathematical statistics. Boston: Academic Press.
- Pratt, J. W. and R. Schlaifer (1984, March). On the nature and discovery of structure. *Journal of the American Statistical Association* 79(385), 9–33.
- Quandt, R. E. (1958, December). The estimation of the parameters of a linear regression system obeying two separate regimes. *Journal of the American Statistical Association* 53(284), 873–880.
- Quandt, R. E. (1988). *The Econometrics of Disequilibrium*. New York: Blackwell.

- Robins, J. (1986). A new approach to causal inference in mortality studies with a sustained exposure period: Application to control of the healthy worker survivor effect. *Mathematical Modelling* 7(9–12), 1393–1512.
- Rosen, S. (1986). The theory of equalizing differences. In O. Ashenfelter and R. Layard (Eds.), *Handbook of Labor Economics*, Volume 1, pp. 641–692. New York: North-Holland.
- Rosenbaum, P. R. and D. B. Rubin (1983, April). The central role of the propensity score in observational studies for causal effects. *Biometrika* 70(1), 41–55.
- Roy, A. (1951, June). Some thoughts on the distribution of earnings. *Oxford Economic Papers* 3(2), 135–146.
- Rubin, D. B. (1974, October). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology* 66(5), 688–701.
- Rubin, D. B. (1978, January). Bayesian inference for causal effects: The role of randomization. *Annals of Statistics* 6(1), 34–58.
- Schennach, S. M. (2020). Mismeasured and unobserved variables. In S. N. Durlauf, L. P. Hansen, J. J. Heckman, and R. L. Matzkin (Eds.), *Handbook of Econometrics, Volume 7A*, Volume 7 of *Handbook of Econometrics*, pp. 487–565. Elsevier.
- Shpitser, I. and J. Pearl (2006, November). Identification of joint interventional distributions in recursive semi-markovian causal models. In *Proceedings of the 21st National Conference on Artificial Intelligence and the 18th Innovative Applications of Artificial Intelligence Conference, AAAI-06/IAAI-06*, Proceedings of the National Conference on Artificial Intelligence, pp. 1219–1226. 21st National Conference on Artificial Intelligence and the 18th Innovative Applications of Artificial Intelligence Conference, AAAI-06/IAAI-06 ; Conference date: 16-07-2006 Through 20-07-2006.
- Shpitser, I. and J. Pearl (2009). Effects of treatment on the treated: Identification and generalization. In *Proceedings of the 25th Conference on Uncertainty in Artificial Intelligence, UAI 2009*, Proceedings of the 25th Conference on Uncertainty in Artificial Intelligence, UAI 2009, pp. 514–521. AUAI Press.
- Tamer, E. (2003, January). Incomplete simultaneous discrete response model with multiple equilibria. *Review of Economic Studies* 70(1), 147–165.
- Telser, L. G. (1964, September). Iterative estimation of a set of linear regression equations. *Journal of the American Statistical Association* 59(307), 845–862.
- Theil, H. (1953). *Estimation and Simultaneous Correlation in Complete Equation Systems*. The Hague: Central Planning Bureau. Mimeographed memorandum.
- Theil, H. (1958). *Economic Forecasts and Policy*. Number 15 in Contributions to Economic Analysis. Amsterdam: North-Holland Publishing Company.
- Theil, H. (1971). *Principles of Econometrics*. New York: Wiley.

- Tinbergen, J. (1930, October). Bestimmung und deutung von angebotskurven ein beispiel. *Zeitschrift für Nationalökonomie* 1(5), 669–679.
- Tinbergen, J. (1939, January). *Statistical Testing of Business Cycle Theories: Part II: Business Cycles in the United States of America, 1919–1932*. Geneva: League of Nations, Economic Intelligence Service.
- Vytlačil, E. J. (2002, January). Independence, monotonicity, and latent index models: An equivalence result. *Econometrica* 70(1), 331–341.
- Wright, S. (1921). Correlation and causation. *Journal of Agricultural Research* 20, 557–585.
- Wright, S. (1934). The method of path coefficients. *Annals of Mathematical Statistics* 5(3), 161–215.
- Yamamoto, T. (2014). Identification and estimation of causal mediation effects with treatment noncompliance. *Unpublished Manuscript, MIT Department of Political Science*.