# Causal Analysis After Haavelmo

**James Heckman** and
The University of Chicago, Department of Economics, 1126 E. 59th St. Chicago, IL 60637 (773) 702-0634

**Rodrigo Pinto**[*]
The University of Chicago, Department of Economics, 1155 E. 60th St., Room 215, Chicago, IL 60637 (773) 702-9436

James Heckman: jjh@uchicago.edu; Rodrigo Pinto: rodrig@uchicago.com

## Abstract

Haavelmo's seminal 1943 and 1944 papers are the first rigorous treatment of causality. In them, he distinguished the definition of causal parameters from their identification. He showed that causal parameters are defined using *hypothetical* models that assign variation to some of the inputs determining outcomes while holding all other inputs fixed. He thus formalized and made operational Marshall's (1890) *ceteris paribus* analysis. We embed Haavelmo's framework into the recursive framework of Directed Acyclic Graphs (DAGs) used in one influential recent approach to causality (Pearl, 2000) and in the related literature on Bayesian nets (Lauritzen, 1996). We compare the simplicity of an analysis of causality based on Haavelmo's methodology with the complex and nonintuitive approach used in the causal literature of DAGs—the "*do-calculus*" of Pearl (2009). We discuss the severe limitations of DAGs and in particular of the *do-calculus* of Pearl in securing identification of economic models. We extend our framework to consider models for simultaneous causality, a central contribution of Haavelmo. In general cases, DAGs cannot be used to analyze models for simultaneous causality, but Haavelmo's approach naturally generalizes to cover them.

## 1 Trygve Haavelmo's Causality

Trygve Haavelmo made fundamental contributions to understanding the formulation and identification of causal models. In two seminal papers (1943 (1944), he formalized the distinction between correlation and causation,[1] laid the foundation for counterfactual policy

[17]We use the term "completeness" in the sense of Koopmans et al. (1950); i.e., the existence of a local solution of Equations (7a) and (7b). This concept is to be distinguished from the notion of completeness in the nonparametric IV literature (Newey and Powell, 2003) or in hypothesis testing (Lehmann and Romano, 2005).

analysis and distinguished the concept of "fixing" from the statistical operation of conditioning—a central tenet of structural econometrics. He developed an empirically operational version of Marshall's notion of *ceteris paribus* (1890), which is a central notion of economic theory, even though Haavelmo never explicitly used that terminology.

In Haavelmo's framework, the causal effects of inputs on outputs are determined by the impacts of *hypothetical* manipulations of inputs on outputs which he distinguishes from correlations between inputs and outputs in observational data. The causal effect of an input is defined using a hypothetical model that abstracts from the empirical data generating process by making hypothetical variation in inputs that are independent of all other determinants of outputs. As a consequence, Haavelmo's notion of causality relies on a thought experiment in which the model that governs the observed data is extended to allow for independent manipulation of inputs, irrespective of whether or not they vary independently in the data.

Haavelmo formalized Frisch's notion that "causality is in the mind."[2] Causal effects are not empirical statements or descriptions of actual worlds, but descriptions of hypothetical worlds obtained by varying—hypothetically—the inputs determining outcomes. Causal relationships are often suggested by observed phenomena, but they are abstractions from it.[3]

This paper revisits Haavelmo's notions of causality using the mathematical language of Directed Acyclic Graphs (DAGs). We start with a recursive framework less general than that of Haavelmo (1943). This allows us to represent causal models as Directed Acyclic Graphs which are intensively studied in the literature on Bayesian networks (Howard and

---

[1]To our knowledge, the first recorded statement of the distinction between correlation and causation is due to Fechner (1851), who distinguished "causal dependency" from what he called "functional relationship." See Heidelberger (2004, p. 102). In later work, Yule (1895, footnote 2, p. 605) discussed the distinction between correlation and causation in a discussion of the effect of relief payments on pauperism. We thank, respectively, Olav Bjerkholt and Steve Stigler for these references.

[2]This notion is central to structural econometrics. It was developed by Frisch and participants in his laboratory, going back to at least 1930:

> "… 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, an intellectual trick … which has proved itself to be a useful weapon … the 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, published 2011)

Writing in the heyday of the Frisch-Haavelmo-inspired Cowles Commission in the late 1940s, Koopmans and Reiersøl distinguished descriptive statistical inference from structural estimation in the following statement.

> "*In many fields the objective of the investigator's inquisitiveness is not just a "population" in the sense of a distribution of observable variables, but a physical structure projected behind this distribution, by which the latter is thought to be generated. The word "physical" is used merely to convey that the structure concept is based on the investigator's ideas as to the "explanation" or "formation" of the phenomena studied, briefly, on his theory of these phenomena, whether they are classified as physical in the literal sense, biological, psychological, sociological, economic or otherwise.*" (Koopmans and Reiersøl 1950, p. 165)

See Simon (1953), Heckman (2008) and Freedman (2010), for later statements of this point of view.

[3]All models—empirical or hypothetical—are idealized thought experiments. There are no formalized rules for creating models, causal or empirical. Analysts may differ about the inputs and relationships in either type of model. A model is more plausible the more phenomena it predicts and the deeper are its foundations in established theory. Causal models are idealizations of empirical models which are in turn idealizations of phenomena. Some statisticians reject the validity of hypothetical models and seek to define causality using empirical methods (Sobel, 2005). As an example we can cite the "Rubin model" of Holland (1986), which equates establishing causality with the empirical feasibility of conducting experiments. This approach confuses the definition of causal parameters with their identification from data. See Heckman (2005, 2008) for a discussion of this approach.

Matheson, 1981; Lauritzen, 1996; Pearl, 2000). We then consider the general non-recursive framework of Haavelmo (1943, 1944) which cannot, in general, be framed as DAGs.

Following Haavelmo, we define hypothetical models that are used to generate causal parameters as idealizations of empirical models that govern the data generating processes. This facilitates discussion of causal concepts such as "fixing" using an intuitive approach that draws on Haavelmo's notion of causality. Identification relies on linking the parameters defined in a hypothetical model using data generated by an empirical model.

This paper makes the following contributions to the literature on causality: (1) We build a framework for the study of causality inspired by Haavelmo's concept of hypothetical variation of inputs. (2) In doing so, we express Haavelmo's notion of causality in the mathematical language of DAGs. (3) For this class of models, we compare the simplicity of Haavelmo's framework with the cumbersome and nonintuitive causal framework for the *do-calculus* proposed by Pearl (2000), which is beginning to be used in economics (see, e.g., Margolis et al., 2012; White and Chalak, 2009). (4) We discuss the limitations of the use of DAGs for econometric identification. We show that even in recursive models, the methods that rely solely on the information in DAGs do not exploit identification strategies based on functional restrictions and exclusion restrictions that are generated by economic theory. This limitation produces apparent non-identification in classically identified econometric models. We show how Haavelmo's approach naturally extends to notions of simultaneous causality while the DAG approach is fundamentally recursive.

Our paper is on the methodology of causality. We do not create a new concept of causality, but rather propose a new framework within which to discuss it. We show that Haavelmo's approach is a complete framework for the study of causality that accommodates the main tools of identification used in the current literature in econometrics, whereas an approach based on DAGs does not.

We show that the causal operation of fixing described in Haavelmo (1943) and Heckman (2005, 2008) is equivalent to statistical conditioning when embedded in a hypothetical model that assigns independent variation to inputs with regard to all variables not caused by those inputs. Pearl (2009) uses the term *do* for the concept of fixing a variable. We show the relationship between statistical conditioning in a hypothetical model and the do-operator. Fixing, in our framework, differs from the operation of the do-operator because it targets specific causal links instead of variables that operate across multiple causal links. A benefit of targeting causal links is that it simplifies the analysis of the subsets of causal relationships associated with an input variable when compared to the do-operator. An analysis of causality based on Haavelmo's approach uses standard rules of probability to define and analyze causal parameters. In contrast, the complex and nonintuitive analysis of Pearl invents new and cumbersome nonstandard rules for analyzing causal models that do not cover a variety of models identified under the Haavelmo approach.

Haavelmo's approach allows for a precise yet intuitive definition of causal effects. With it, analysts can identify causal effects by applying standard statistical tools. In contrast with the do-calculus, application of Haavelmo's concepts eliminates the need for additional extra-

statistical graphical/statistical rules to achieve identification of causal parameters. Haavelmo's approach also covers the case of simultaneous causality in its full generality whereas frameworks for causal analysis currently used in statistics cannot, except through introduction and application of *ad hoc* rules.

This paper is organized in the following way. Section 2 reviews Haavelmo's causal framework. Section 3 uses a modern framework of causality to assess Haavelmo's contributions to the literature. Section 4 examines how application of this framework differs from Pearl's do-calculus (2009) and enables analysts to apply the standard tools of probability and statistics without having to invent new extra-statistical rules. It gives an example of the identification of causal effects that considers Pearl's "Front-Door" criteria and the nonidentifiability of the instrumental variables model using the rules of the do-calculus. Section 5 discuss the limitations of DAGs in implementing the variety of sources of identification available to economists. We focus on the simplest cases of confounding models where instrumental variables are available. Section 6 extends the discussion to a simultaneous equations framework. Section 7 concludes.

## 2 Haavelmo's Causal Framework

We review the key concepts of causality developed by Haavelmo (1943, 1944)—starting with a recursive model. A causal model is based on a system of structural equations that define causal relationships among a set of variables. In the language of Frisch (1938), these structural equations are *autonomous* mechanisms represented by deterministic functions mapping inputs to outputs. By autonomy we mean, as did Frisch, that these relationships remain invariant under external manipulations of their arguments. They are functions in the ordinary usage of the term in mathematics. They produce the same values of the outcomes when inputs are assigned to a fixed set of values, however those values are determined. Even though the functional form of a structural equation may be unknown, the causal directions among the variables of a structural equation are assumed to be known. They are determined by thought experiments that may sometimes be validated in data. The variables chosen as arguments in a structural equation are assumed to account for all causes of the associated output variable.

Haavelmo developed his work on causality for aggregate economic models. He considered mean causal effects and, for the sake of simplicity, invoked linearity, assumed uniformity of responses to inputs across agents, and focused on continuous variables. More recent approaches generalize his framework.

Haavelmo formalized the distinction between correlation and causation using a simple model. In order to examine his ideas, consider three variables *Y, X, U* associated with error terms $\boldsymbol{\varepsilon} = (\varepsilon_U, \varepsilon_X, \varepsilon_Y)$ such that *X, Y* are observed by the analyst while variables *U, ε* are not. [4] He assumed that *U* is a confounding variable that causes *Y* and *X*. We represent this model through the following structural equations:

---

[4]This framework allows for uncertainty on the part of agents if realizations of the uncertain variables are captured through variables *X* and *U*. In that sense the model can be characterized as a method for examining *ex-post* relationships between variables. For a discussion of causal analysis of *ex-post* versus *ex-ante* models, see, e.g., Hansen and Sargent (1980) and Heckman (2008).

$$Y = f_Y(X, U, \varepsilon_Y), X = f_X(U, \varepsilon_X), \qquad \text{and } U = f_U(\varepsilon_U) \qquad ,$$

where $\varepsilon$ is a vector of mutually independent error terms with cumulative distribution function $Q_\varepsilon$. Thus, if $X$, $U$, $\varepsilon_Y$ take values of $x$, $u$, $e_Y$, then $Y$ must take the value $y = f_Y(x, u, e_Y)$. By iterated substitution we can express all variables in terms of $\varepsilon$. Moreover, the mutual independence assumption of error terms implies that $\varepsilon_Y$ is independent of $(X, U)$ as $X = f_X(f_U(\varepsilon_U), \varepsilon_X)$ and $U = f_U(\varepsilon_U)$. Notationally, we write $(X, U) \perp\!\!\!\perp \varepsilon_Y$, where $\perp\!\!\!\perp$ denotes statistical independence. In the same fashion, we have that $\varepsilon_X \perp\!\!\!\perp U$ but $X$ is not independent of $\varepsilon_U$.

Haavelmo defines the causal effect of $X$ on $Y$ as being generated by a *hypothetical manipulation* of variable $X$ that does not affect the values that $U$ or $\varepsilon$ take. This is called *fixing $X$* by a hypothetical manipulation.[5] Notationally, outcome $Y$ when $X$ is fixed at $x$ is denoted by $Y(x) = f_Y(x, U, \varepsilon_Y)$ and its expectation is given by $\mathbf{E}_{(U, \varepsilon_Y)}(Y(x)) = \mathbf{E}(f(x, U, \varepsilon_Y))$, where $\mathbf{E}_{(U, \varepsilon_Y)}(\cdot)$ means expectation over the distribution of random variables $U$ and $\varepsilon_Y$. The average causal effect of $X$ on $Y$ when $X$ takes values $x$ and $x'$ is given by $\mathbf{E}_{(U, \varepsilon_Y)}(Y(x)) - \mathbf{E}_{(U, \varepsilon_Y)}(Y(x'))$. For notational simplicity, we henceforth suppress the subscript on $\mathbf{E}$ denoting the random variable with respect to which the expectation is computed.

Conditioning is a statistical operation that accounts for the dependence structure in the data. Fixing is an abstract operation that assigns independent variation to the variable being "fixed." The standard linear regression framework is convenient for illustrating these ideas and in fact is the one used by Haavelmo (1943).

Consider the standard linear model $Y = X\beta + U + \varepsilon_y$ where $\mathbf{E}(\varepsilon_Y) = 0$ represent the data generating process for $Y$. The expectation of outcome $Y$ when $X$ is *fixed* at $x$ is given by $\mathbf{E}(Y(x)) = x\beta + \mathbf{E}(U)$. This equation corresponds to Haavelmo's (1943) hypothetical model. The expectation of $Y$ when $X$ is *conditioned* on $x$ is given by $\mathbf{E}(Y|X=x) = x\beta + \mathbf{E}(U|X=x)$, as $\mathbf{E}(\varepsilon_Y|X=x) = 0$ because $\varepsilon_Y \perp\!\!\!\perp X$. If $\mathbf{E}(U|X=x) = 0$ and elements of $X$ are not collinear, then OLS identifies $\beta$ and $\mathbf{E}(Y|X=x) = \mathbf{E}(Y(x)) = x\beta$ and $\beta$ generates a causal parameter: the average treatment effect of a change in $X$ on $Y$. Specifically, $(x - x')\beta$ is the average difference between the expectation of $Y$ when $X$ is fixed at $x$ and $x'$.

The difficulty of identifying the average causal effect of $X$ on $Y$ when $\mathbf{E}(U|X) \neq 0$ (and thereby $\mathbf{E}(Y|X=x) \neq \mathbf{E}(Y|(x))$) stems from the potential confounding effects of unobserved variable $U$ on $X$. In this case, the standard Least Squares estimator does not generate an autonomous causal or structural parameter because $plim(\hat{\beta}) = \beta + \text{cov}(X, U)/\text{var}(X)$ depends on the covariance between $X$ and $U$. While the concept of a causal effect

---

[5]Haavelmo (1943) did not explicitly use the term "fixing." He set $U$ (in our notation) to a specified value and manipulated $X$ in his "hypothetical model." Specifically Haavelmo set $U = 0$ but the point of evaluation is irrelevant in the linear case he analyzed.

does not rely on the properties of the data generating process, the identification of causal effects does.

Without linearity, one needs an assumption stronger than $\mathbf{E}(U|X=x)=0$ to obtain $\mathbf{E}(Y|X=x)=\mathbf{E}(Y(x))$. Indeed if one assumes no confounding effects of $U$, that is to say that $X$ and $U$ are independent ($X \perp\!\!\!\perp U$), then one can show that fixing is equivalent to statistical conditioning:

$$
\begin{aligned}
\mathbf{E}(Y|X=x) \quad &= \int f_Y(x,u,\varepsilon_Y)\, dQ_{(U,\varepsilon_Y)|X=x}(u,\varepsilon_Y) \\
&= \int f_Y(x,u,\varepsilon_Y)\, dQ_U(u)\, dQ_{\varepsilon_Y}(\varepsilon_Y) \\
&= \mathbf{E}(f_Y(x,U,\varepsilon_Y)) \\
&= \mathbf{E}(Y(x)),
\end{aligned}
$$

where $Q_{(U,\ \varepsilon_Y)|X=x}(u,\ \varepsilon_Y)$ denotes the cumulative joint distribution function of $U$, $\varepsilon_Y$ conditional on $X = x$ and the second equality comes from as the fact that $U$, $X$ and $\varepsilon_Y$ are mutually independent. If $X \perp\!\!\!\perp (U,\ \varepsilon_Y)$ holds, we can use observational data to identify the mean value of $Y$ fixing $X = x$ by evaluating the expected value of $Y$ conditional on $X = x$. Note that in general, the value obtained depends on the functional form of $f_Y(x,\ u,\ \varepsilon_Y)$.

Haavelmo's notation has led to some confusion in the statistical literature. His argument was aimed at economists of the 1940s and does not use modern notation. Haavelmo's key definitions and ideas are given by examples rather than by formal definitions. We restate and clarify his framework in this paper.

To simplify the exposition, assume that all variables are discrete and let *Pr* denote their probability measure. The factorization of the joint distribution of *Y, U* conditional on *X* is given by *Pr(Y, U|X = x) = Pr(Y|U, X = x)Pr(U|X = x)*. In contrast, in the abstract operation of fixing *X* is assumed not to affect the marginal distribution of *U*. That is to say that *U(x) = U*. Therefore the joint distribution of *Y, U* when *X* is fixed at *x* is given by *Pr(Y(x), U(x)) = Pr(Y(x), U) = Pr(Y|U, X = x)Pr(U)*.

Fixing lies outside the scope of standard statistical theory and is often a source of confusion. Indeed, even though the probabilities *Pr(Y|U, X = x)* and *Pr(U)* are well defined, neither the causal operation of fixing nor the resulting joint distribution follow from standard statistical arguments.[6] Conditioning *is* equivalent to fixing under independence of *X* and *U*. In this case the conditional joint distribution of *Y* and *U* becomes *Pr(Y, U|X = x) = Pr(Y|U, X = x)Pr(U|X = x) = Pr(Y|U, X = x)Pr(U)*.

To gain more intuition on the difference between fixing and conditioning, express the conditional expectation $\mathbf{E}(Y|X=x)$ as the integral across $\varepsilon$ over a restricted set $\mathscr{A}^C$. By iterated substitution, we can write *Y* as $Y = f_Y(f_X(f_U(\varepsilon_U), \varepsilon_X), fu(\varepsilon_u),\ \varepsilon_Y)$. Thus

[6]See Pearl (2009) and Spirtes et al. (2000) for discussions.

$$\mathbf{E}(Y|X=x) = \frac{\int_{\mathscr{A}^C} f_Y(f_X(f_U(\varepsilon_U), \varepsilon_X), f_U(\varepsilon_U), \varepsilon_Y) dQ_\varepsilon(\varepsilon)}{\int_{\mathscr{A}^C} dQ_\varepsilon(\varepsilon)} \quad (1)$$

$$\text{where } \mathscr{A}^C = \{\varepsilon = (\varepsilon_U, \varepsilon_X, \varepsilon_Y) \in \text{supp}(\varepsilon); f_X(f_U(\varepsilon_U), \varepsilon_X) = x\}. \quad (2)$$

Fixing, on the other hand, is written as the integral across $\varepsilon$ over its full support:

$$\mathbf{E}(Y(x)) = \frac{\int_{\mathscr{A}^F} f_Y(x, f_U(\varepsilon_U), \varepsilon_Y) dQ_\varepsilon(\varepsilon)}{\int_{\mathscr{A}^F} dQ_\varepsilon(\varepsilon)} \quad (3)$$

$$\text{where } \mathscr{A}^F = \{\varepsilon = (\varepsilon_U, \varepsilon_X, \varepsilon_Y) \in \text{supp}(\varepsilon)\} \text{ and } \int_{\mathscr{A}^F} dQ_\varepsilon(\varepsilon) = 1. \quad (4)$$

Fixing differs from conditioning in terms of the difference in the integration sets $\mathscr{A}^F$ and $\mathscr{A}^C$. While conditional expectation (1) is a standard operation in statistics, the operation used to define fixing is not. Equation (1) is an expectation conditional on the event $f_X(f_U(\varepsilon_U), \varepsilon_X) = x$, which affects the integration set $\mathscr{A}^C$ given in (2). Fixing (3), on the other hand, integrates the function $f_Y(x, f_U(\varepsilon_U), \varepsilon_Y)$ across the whole support of $\varepsilon$ given in (4). The inconsistency between fixing and conditioning in the general case comes from the fact that fixing $X$ is equivalent to setting the expression $f_X(f_U(\varepsilon_U), \varepsilon_X)$ to $x$ without changing the probability measures of $\varepsilon_U$, $\varepsilon_X$ associated with the operation of conditioning on the event $X = x$.

This paper interprets Haavelmo's approach by introducing a hypothetical model that enables analysts to examine fixing using standard tools of probability. The *hypothetical model* departs from the data generating process by exploiting autonomy and creating a *hypothetical* variable that has the desired property of independent variation with regard to $U$. The hypothetical model is an idealization of the empirical model. Standard statistical tools apply to both the data generating process and the hypothetical model.

To formalize Haavelmo's notions of causality, let a hypothetical model with error terms $\varepsilon$ and four variables including $Y$, $X$, $U$ but also a new variable $\tilde{X}$ with the property that $\tilde{X} \perp\!\!\!\perp (X, U, \varepsilon)$.[7] Invoking autonomy, the hypothetical model shares the same structural equation as the empirical one but departs from it by replacing $X$ with an $\tilde{X}$-input, namely $Y = f_Y(\tilde{X}, U, \varepsilon_Y)$. The hypothetical model is not a wildly speculative departure from the empirical data generating process but an expanded version of it. Thus $(Y|X = x, U = u) = f_Y(x, u, \varepsilon_Y)$ in the empirical model and $(Y|\tilde{X} = x, U = u) = f_Y(x, u, \varepsilon_Y)$ in the hypothetical model. The hypothetical model has the same marginal distribution of $U$ as the empirical model. The joint distributions of variables in the empirical model $Pr_E$ and the hypothetical model $Pr_H$ may differ.

---

[7]We could express $\tilde{X} = f_X(\tilde{\varepsilon}_X)$ to be notationally consistent.

The hypothetical model clarifies the notion of fixing in the empirical model. Fixing in the empirical model is based on non-standard statistical operations. However, the distribution of the outcome $Y$ when $X$ is fixed at $x$ in the empirical model can be interpreted as standard statistical conditioning in the hypothetical model, namely, $Pr_{\mathrm{E}}(Y(x)) = Pr_{\mathrm{H}}(Y|\tilde{X} = x)$. The next section formalizes these ideas using one modern language of causality.[8]

## 3 Recasting Haavelmo's Ideas

We recast Haavelmo's model in the framework of Directed Acyclic Graphs (DAGs). DAGs are studied in Bayesian Networks (Howard and Matheson, 1981; Lauritzen, 1996) and are often used to define and estimate causal relationships (Lauritzen, 2001). The literature on causality based on DAGs was advanced by Judea Pearl (2000, 2009).[9]

In this fundamentally recursive framework, a causal model consists of a set of variables $\mathscr{T} = \{T_1, \ldots, T_n\}$ associated with a set of mutually independent error terms $\varepsilon = \{\varepsilon_1, \ldots, \varepsilon_n\}$ and a system of autonomous structural equations $\{f_1, \ldots, f_n\}$. Variable set $\mathscr{T}$ includes both observed and unobserved variables. Variable set $\mathscr{T}$ also include both external and internal variables. We clarify these concepts in the following way.

Causal relationships between a dependent variable $T_i \in \mathscr{T}$ and its arguments are defined by $T_i = f_i(Pa(T_i), \varepsilon_i)$, where $Pa(T_i) \subset \mathscr{T}$ and $\varepsilon_i \in \varepsilon$ are called parents of $T_i$ and are said to directly cause $T_i$. If $Pa(T) = \varnothing$ then variable $T$ is not caused by any variable in $\mathscr{T}$. In this case, $T$ is an *external variable* determined outside the system, otherwise the variable is called an *internal or endogenous variable*. The error terms in $\varepsilon$ are not caused by any variable and are introduced to avoid degenerate conditioning statements among variables in $\mathscr{T}$. For simplicity of notation, we keep the error terms $\varepsilon$ implicit, except when it clarifies matters to do so. We assume that all random variables in this section and the next are discrete valued although this requirement is easily relaxed.

Causal relationships are represented by a graph $G$ where each node corresponds to a variable $T \in \mathscr{T}$. Nodes are connected by arrows from $Pa(T)$ to $T$ and represent causal influences among variables. Descendants of a variable $T$, i.e., $D(T) \subset \mathscr{T}$, consist of all variables connected to $T$ by arrows of the same direction arising from $T$. Graph $G$ is called a DAG if no variable is a descendant of itself, i.e., $T \notin D(T), \forall\, T \in \mathscr{T}$. Observe that this assumption rules out simultaneity—a central feature of Haavelmo's approach. Children of a variable $T$ are the set of variables that have $T$ as a parent, namely, $Ch(T) = \{T' \in \mathscr{T}; T \in Pa(T')\}$.

Causal relationships are translated into statistical relationships in a DAG through a property termed the Local Markov Condition (LMC) (Kiiveri et al., 1984; Lauritzen, 1996). LMC states that a variable is independent of its non-descendants conditional on its parents. LMC (5) also holds among variables in $\mathscr{T}$ under the assumption that error terms $\{\varepsilon_1, \ldots, \varepsilon_n\}$ are mutually independent (Pearl, 1988; Pearl and Verma, 1994), namely:

---

[8]Frisch's (1938) notion of invariance used by Haavelmo is called SUTVA in one model of causality popular in statistics. See Holland (1986) and Rubin (1986).
[9]Chalak and White (2012) present generalizations of this approach.

$$\mathbf{LMC}:\text{for all } \mathbf{T} \in \mathscr{T}, \mathbf{T} \perp\!\!\!\perp (\mathscr{T} \backslash \mathbf{D}(\mathbf{T}) \cup T) | Pa(T). \quad (5)$$

We use Dawid's (1979) notation to denote conditional independence. If $W$, $K$, $Z$ are subsets of $\mathscr{T}$, the expression $W \perp\!\!\!\perp K|Z$ means that each variable in $W$ is statistically independent of each variable in $K$ conditional on all variables in $Z$. The conditional independence relationships generated by LMC (5) can be further manipulated using the Graphoid relations. [10] An important use of LMC (5) is to factorize the joint distribution of variables $Pr(T_1, \ldots, T_n)$. Under a recursive model, we can assume without loss of generality that variables $(T_1, \ldots, T_n, \ldots, T_N)$ are ordered so that $(T_1, \ldots, T_{n-1})$ are non-descendants of $T_n$ and thereby $Pa(T_n) \subset (T_1, \ldots, T_{n-1})$. Thus,

$$\Pr(T_1, \ldots, T_n) = \prod_{T_n \in \mathscr{T}} \Pr(T_n | T_1, \ldots, T_{n-1}) = \prod_{T_n \in \mathscr{T}} \Pr(T_n | Pa(T_n)), \quad (6)$$

where the last equality comes from applying LMC (5).

Table 1 uses the Haavelmo model described in Section 2 to illustrate the concepts discussed here. Table 1 presents two models and six panels separated by a series of horizontal lines. The first panel names the models. The second panel presents the structural equations generating the models. Columns 1 and 2 are based on structural equations that have the same functional form, but different inputs. The third panel represents the associated model as a DAG. Squares represent observed variables, circles represent unobserved variables. (Except in the first panel, the components of $\varepsilon$ are kept implicit in the table.) The fourth panel displays the parents in $\mathscr{T}$ for each variable. The fifth panel shows the conditional independence relationships generated by the application of LMC (5), and the sixth and final panel presents the factorization of the joint distribution.

We use the framework presented above to discuss the concept of fixing in greater generality. According to Section 2, we define the causal operation of fixing a variable in a model represented by a graph $G$ by the intervention that sets a value to this variable in $\mathscr{T}$ in a fashion that does not affect the distribution of its non-descendants. In other words, fixing a random variable (or a set of random variables) $X \in \mathscr{T}$ to $x$ translates to setting $X = x$ for *all* $X$-inputs in the structural equations associated with variables in $Ch(X)$. Pearl (2009) uses the term *doing* for what we call *fixing*. We use his notation in writing Equation (7), below. The post-intervention distribution of variables in $\mathscr{T}$ when $X$ is fixed at $x$ is given by

---

[10]The Graphoid relationships are a set of elementary conditional independence relationships presented by Dawid (1979):

| | |
|---|---|
| Symmetry | $X \perp\!\!\!\perp Y|Z \Rightarrow Y \perp\!\!\!\perp X|Z.$ |
| Decomposition: | $X \perp\!\!\!\perp (W, Y)|Z \Rightarrow X \perp\!\!\!\perp Y|Z.$ |
| Weak Union: | $X \perp\!\!\!\perp (W, Y)|Z \Rightarrow X \perp\!\!\!\perp W|(Y, Z).$ |
| Contraction: | $X \perp\!\!\!\perp Y|Z$ and $X \perp\!\!\!\perp W|(Y, Z) \Rightarrow X \perp\!\!\!\perp (W, Y)|Z.$ |
| Intersection: | $X \perp\!\!\!\perp W|(Y, Z)$ and $X \perp\!\!\!\perp Y|(W, Z) \Rightarrow X \perp\!\!\!\perp (W, Y)|Z.$ |
| Redundancy: | $X \perp\!\!\!\perp Y|X.$ |

The intersection relation is only valid for variables with strictly positive probability distributions. See also Dawid (2001).

$$\Pr(\mathscr{T}\backslash\{X\}|do(X)=x)=\prod_{T\in\mathscr{T}\backslash\{\{X\}\cup Ch(X)\}}\Pr(T|Pa(T))\prod_{T\in Ch(X)}\Pr(T|Pa(T)\backslash\{X\},X=x). \quad (7)$$

Versions of Equation (7) can be found in Pearl (2001), Spirtes et al. (2000), and Robins (1986). In this instance, $do(X) = x$ is equivalent to conditioning $X$ at $\tilde{X} = x$.

As noted in Section 2, standard arguments based on statistical conditioning are unable to describe the probability laws governing the fixing operation used in Equation (7). Our solution to this problem draws on Haavelmo's insight that causality is a property of hypothetical models in which causal effects on output variables are generated through hypothetical independent variations of inputs. Specifically, we are able to map causal manipulations of the fixing operation into standard statistical language by formalizing the concept of a hypothetical model in Section 3.1.

## 3.1 The Hypothetical Model

We formalize the concept of a hypothetical model and study its properties. The notions discussed here constitute our theoretical basis to examine causal effects. To recall, we use the term *empirical model* to designate the data generating process and the term *hypothetical model* to designate the model used to characterize causal effects.

The hypothetical model is generated from an empirical model. It shares the same structural equations and same distribution of error terms as the empirical model. The hypothetical model differs from the empirical model in two ways. First, it appends to the empirical model an external variable (or a set of external variables) termed a hypothetical variable(s). Second, it replaces the action of existing inputs. If $X \in \mathscr{T}$ is the target variable to be fixed in the empirical model, then the newly created hypothetical variable $\tilde{X}$ replaces the $X$-input of one, some or all variables in $Ch(X)$. In other words, children of $X$ in the empirical model will have their $X$-input replaced by a $\tilde{X}$-input in the hypothetical model. We assume that $X$ and $\tilde{X}$ have common supports.

Table 1 illustrates the concept of a hypothetical model using the Haavelmo model introduced in Section 2. Column 1 presents the Haavelmo empirical model while Column 2 presents its associated hypothetical model.

For the sake of clarity, we use $G_E$ for the DAG representing the empirical model and $\mathscr{T}_E$ for its associated set of variables. We use $Pa_E$, $D_E$, $Ch_E$ for the parents, descendants, and children with DAG $G_E$. We use $Pr_E$ for the probability measure of variables in $\mathscr{T}_E$. For the corresponding counterparts in the hypothetical model we use $G_H$, $\mathscr{T}_H$, $Pa_H$, $D_H$, $Ch_H$, and $Pr_H$.

We now list some salient features of the hypothetical model. Let $\tilde{X}$ denote the hypothetical variable (or variables) associated with $X \in \mathscr{T}_E$. We expand the list of variables in the hypothetical model so that $\mathscr{T}_H = \mathscr{T}_E \cup \{\tilde{X}\}$. The hypothetical variable can replace some or all of the input $X$ for variables in $Ch_E(X)$, i.e., $Ch_H(X) \subseteq Ch_E(X)$. Children of $X$ in the empirical model can be partitioned among $X$ and $\tilde{X}$ in the hypothetical model: $Ch_E(X) =$

$Ch_H(X) \cup Ch_H(X)\tilde{}.$[11] As a consequence we also have that $D_E(X) = D_H(X) \cup D_H(X)\tilde{},$ that is, $X$-descendants of the empirical model constitute the $X$ and $\tilde{X}$ descendants in the hypothetical model. Parental sets of the hypothetical model are defined by $Pa_H(T) = Pa_E(T) \ \forall \ T \in \mathcal{T}_E \backslash Ch_H(X)\tilde{}$ and $Pa_H(T) = \{Pa_E(T)\backslash\{X\}\}\cup\{X\} \ \forall \ T \in Ch_H(X)\tilde{}.$ Moreover, $\tilde{X}$ is an external variable, that is, $Pa_H(\tilde{X}) = \varnothing.$ The hypothetical model is also a DAG. Thus LMC (5) holds and the joint distribution of the variables in $\mathcal{T}_H$ can be factorized using equation (6). By sharing the same structural equations and distribution of error terms, the conditional probabilities of the hypothetical model can be written as

$$\Pr_H(T|Pa_H(T)) = \Pr_E(T|Pa_E(T)) \forall T \in \boldsymbol{\mathcal{T}}_E \backslash Ch_H(\tilde{X}) \quad (8)$$

and

$$\Pr_H(T|Pa_H(T)\backslash\{\tilde{X}\}, \tilde{X}=x) = \Pr_E(T|Pa_E(T)\backslash\{X\}, X=x) \forall T \in Ch_H(\tilde{X}). \quad (9)$$

Equations (8)–(9) arise because the distribution of a variable $T \in \mathcal{T}_E$ conditional on its parents is determined by the distribution of its error terms, which is the same for hypothetical and empirical models.

We now link the probability measures of the empirical and hypothetical models. Theorem T-1 uses LMC (5) and Equation (8) to show that the distribution of non-descendants of $\tilde{X}$ are the same in both hypothetical and empirical models:

**Theorem T-1.** Let $\tilde{X}$ be the hypothetical variable in the hypothetical model represented by $G_H$ associated with variable $X$ in empirical model $G_E$. Let $W, Z$ be any disjoint set of variables in $\mathcal{T}_E \backslash D_H(\tilde{X}).$ Then

$$\Pr_H(W|Z) = \Pr_H(W|Z, \tilde{X}) = \Pr_E(W|Z) \forall \{W, Z\} \subset T_E \backslash D_H(\tilde{X}).$$

*Proof.* See Appendix.

Theorem T-1 also holds for the set of variables that are non-descendants of $X$ according to the empirical model, which are a subset of $\mathcal{T}_E \backslash D_H(\tilde{X}).$ Thus, $Pr_H(W|Z) = Pr_H(W|Z, \tilde{X}) = Pr_E(W|Z)$ for all $\{W, Z\} \subset \mathcal{T}_E \backslash D_E(X).$

The following theorem uses Theorem T-1 and Equations (8)–(9) to show that the distribution of variables conditional on $X$ and $\tilde{X}$ taking the same value $x$ in the hypothetical model is equal to the distribution of the variables conditional on $X = x$ in the empirical model:

---

[11] As an example, let a simple empirical model for mediation analysis consist of three variables: an input variable $X$, a mediation variable $M$ caused by $X$ and an outcome of interest $Y$ caused by $X$ and $M$. This model is represented as a DAG in Model 1 of Table 2 and $Ch_E(X) = \{M, Y\}.$ Suppose we are interested in the indirect effect, that is the effect of $X$ on $Y$ that operates exclusively by changes in $M$ while holding the distribution of $X$ unaltered. The hypothetical model for the evaluation of the indirect causal effect assigns the causal link of $X$ on $M$ to the hypothetical variable $\tilde{X}$. Namely $X$ still causes $Y$, but $\tilde{X}$ causes $M$. This hypothetical model is represented by Model 3 of Table 2. In this model $Ch_H(X) = \{Y\}, Ch_H(\tilde{X}) = \{M\}$ and $Ch_E(X) = Ch_H(X) \cup Ch_H(\tilde{X}).$

**Theorem T-2.** Let $\tilde{X}$ be the hypothetical variable in the hypothetical model represented by $G_H$ associated with variable $X$ in empirical model $G_E$ and let $W, Z$ be any disjoint[12] set of variables in $\mathscr{T}_E$. Then

$$\mathrm{Pr}_H(W|Z, X{=}x, \tilde{X}{=}x){=}\mathrm{Pr}_E(W|Z, X{=}x)\forall\{W, Z\} \subset \boldsymbol{\mathscr{T}}_E.$$

*Proof.* See Appendix.[13]

A useful corollary of Theorem T-2 is the method of *matching:*

**Corollary C-2. Matching:** Let $Z$, $W$ be any disjoint set of variables in $\mathscr{T}_E$ and let $\tilde{X}$ be a hypothetical variable in model $G_H$ associated with $X \in \mathscr{T}_E$ in model $G_E$ such that, in the hypothetical model, $X \perp\!\!\!\perp W | (Z, \tilde{X})$, then

$$\mathrm{Pr}_H(W|Z, \tilde{X}{=}x)\mathrm{Pr}_E(W|Z, X{=}x).$$

Proof. See Mathematical Appendix.

Variables $Z$ of C-2 are called matching variables. In statistical jargon, it is said that matching variables solve the problem of confounding effects between a treatment indicator $X$ and outcome $W$. Matching is commonly used to identify treatment effects in propensity score matching models.[14] In these models, the conditional independence relation of Matching C-2 is assumed to be true. Pearl (1993) describes a graphical test called the "Back-Door" criterion that can be applied to a DAG in order to check if a set of variables satisfy the assumptions of Matching C-2. We turn next to an explicit discussion of the *do-calculus*.

### 3.2 Benefits of the Hypothetical Model

The major benefit of the hypothetical model is to allow us to perform causal operations using standard statistical tools. As previously noted, the fixing operation is poorly defined in statistics. Statistical tools such as LMC (5), Graphoid Axioms or the Law of Iterated Expectation do not apply to the fixing operator. The hypothetical model solves this mismatch between causal language and statistical operations because the operation of fixing a variable in the empirical model is easily translated into statistical conditioning in the hypothetical model. This comes as a consequence of the properties of the hypothetical variable, which is defined to have the desired independent variation to generate causal

---

[12]Disjoint (i.e., distinct) from $X$

[13]We also note the following result:

**Corollary C-1.** Let $\tilde{X}$ be uniformly distributed in the support of $X$ and let $W, Z$ be any disjoint set of variables in $\mathscr{T}_E$ then:

$$\mathrm{Pr}_H(W|Z, X{=}\tilde{X}){=}\mathrm{Pr}_E(W|Z)\forall\{W, Z\} \subset \boldsymbol{\mathscr{T}}_E.$$

*Proof.* See Appendix. We thank an anonymous referee for suggesting this result and its proof.

[14]See, e.g., Rosenbaum and Rubin (1983).

effects. In particular, if we replace the *X*-input by a $\tilde{X}$-input for all children of *X*, as suggested by the operation of fixing or "doing," we have that the distribution of an outcome $Y \in \mathscr{T}_E$ of the empirical model when variable *X* is fixed at *x* (for all its children) is equivalent to the distribution of *Y* conditional on the hypothetical variable $\tilde{X}$ being assigned to value *x*. This is captured by the following theorem:

Theorem T-3. Let $\tilde{X}$ be the hypothetical variable in $G_H$ associated with variable *X* in the empirical model $G_E$, such that $Ch_H(\tilde{X}) = Ch_E(X)$, then:

$$\text{Pr}_H(\mathscr{T}_E \setminus \{X\} | \tilde{X}{=}x) = \text{Pr}_E(\mathscr{T}_E \setminus \{X\} | do(X){=}x).$$

*Proof.* See Appendix.

Theorem T-3 avoids the need for defining new mathematical tools that would be necessary to integrate the fixing operator into statistical language. Section 4 illustrates this point by comparing the identification of causal effects in the "Front-Door" model of Pearl (2009) using his "*do-calculus*" and identification using a hypothetical model, which does not require any additional apparatus outside standard statistical analysis.

Another benefit of the hypothetical model is that it clearly distinguishes the characterization of a causal effect from its identification in data. Causal effects of a variable *X* on and outcome *Y* are characterized within the hypothetical model by the distribution of *Y* conditioned on hypothetical variable $\tilde{X}$. Identification of causal effects, on the other hand, requires analysts to relate the hypothetical and empirical distributions in a fashion that allows the evaluation of causal effects examined in the hypothetical model using data generated by the empirical model. This task relies on the statistical properties that connect both models; that is, Equations (8) and (9), Theorems T-1 and T-2, and Corollary C-2. Section 4 illustrates model identification while discussing the differences between the hypothetical model approach and the do-calculus of Pearl (1995).

The hypothetical model does not suppress the variable we seek to fix, but rather creates a new hypothetical variable that allows us to examine a variety of causal effects. This approach provides a natural framework within which to examine counterfactual outcomes that involve both fixing and conditioning. For example, suppose *X* denotes schooling choice: *X* = 1 college education and *X* = 0 otherwise. The treatment-on-the-treated parameter stands for the average causal effect of college education for the sub-sample of agents that choose to go to college. This parameter is readily defined by

$\mathbf{E}_H(Y | \tilde{X}{=}1, X{=}1) - \mathbf{E}_H(Y | \tilde{X}{=}0, X{=}1)$ in the hypothetical model. For more examples of such parameters, see Heckman and Vytlacil (2007a).

In contrast with the do-operator, the hypothetical model targets causal links, not variables. In other words, the hypothetical model allows analysts to target separate causal relationships of *X*. We can choose subsets of variables in *Ch*(*X*) that will be caused by a hypothetical variable $\tilde{X}$, which in turn replaces some of the *X* inputs. The do-operator, on the other hand, targets *all* causal relationships involving the variable *X*. However it is often that case that the

analysis would like to examine a counterfactual outcomes generated by fixing *X* at different levels for different structural equations that have *X* as input.

For example, consider the standard mediation model in which *X* denoted treatment (1 for treated and 0 for control assignments), *M* is a mediator variable that is caused by X and Y is an outcome of interest that is caused by *X* and *M*. Suppose the analysis is interested in distinguishing the effect of *X* on *Y* that operates through *X* itself and through the mediated effect of *X* on *M*. Say we are interested in the expected value of the counterfactual outcome *Y* generated by fixing its X-input to 1 while using the distribution of *M* that would be generated by fixing *X* at 0. We can interpret this counterfactual as the result of a treatment direct effect and a control mediation effect.

By evoking $do(X) = 1$ we fix its value at 1 in both *Y* and *M* inputs. But the counterfactual outcome *Y* is based on the distribution of *M* that is generated by fixing *X* at 0. To solve this problem, we are required to also fix *M* to its desired level. As a consequence, we need to evoke $do(X) = 1$ and $do(M)$. Specifically, we can characterize the suggested counterfactual outcome in terms of the do-operator by:

$$\int \mathbf{E}_{\mathrm{E}}(Y \,|\, do(X){=}1, do(M){=}m)\, dQ_{M|do(X)=0}(m), \quad (10)$$

$Q_{M|do(X)=0}$ denotes the cumulative distribution of the variable *M* when *X* is fixed at 0, that is $(M|do(X) = 0)$. Under the hypothetical model presented in Model 2 of Table 2, the studied counterfactual outcome is easily written as $\mathbf{E}_{\mathrm{H}}(Y \,|\, \tilde{X}{=}0, X{=}1)$.

## 4 Do-Calculus versus Hypothetical Model

This section compares the identification strategy of the *do-calculus* developed by Pearl (1995) with the causal framework based on a hypothetical model. The hypothetical model approach is conceptually simpler, more tractable and more intuitive.

The do-calculus consists of three graphical and statistical rules that operate on the empirical model. These rules supplement standard statistical theory in order to process the information of a causal models that can be represented by a DAG. The goal of the do-calculus rules is to integrate the statistically ill-defined concept of fixing into a standard statistical framework. Examples of this literature are Huang and Valtorta (2006, 2008) and Tian and Pearl (2002, 2003).

Contrary to the hypothetical model approach, the do-calculus requires special graph-theoretic notation outside common statistical usage. Namely, let *X, Y, Z, W* be arbitrary disjoint sets of variables (nodes) in a causal graph *G*, then the graphical operations of the do-calculus are given by:

1. $G_{\overline{X}}$ denotes a modification of DAG *G* obtained by deleting the arrows pointing to *X*;

2. $G_{\underline{X}}$ denotes the modified DAG obtained by deleting the arrows emerging from *X*;

3.  $G_{\overline{X},\underline{Z}}$ denotes the DAG obtained by deleting arrows pointing to *X* and emerging from *Z*.

The graphical operations items 1,2 and 3 above also apply to nodes *Z(W)*, which denotes the set of *Z*-nodes that are not ancestors of any *W*-node in a DAG *G*. Table 3 gives examples of this notation.

Let *G* be a DAG and let *X, Y, Z, W* be any disjoint sets of variables. The do-calculus rules are:

- **Rule 1:** Insertion/deletion of variables:

  $Y \perp\!\!\!\perp Z|(X, W)$ under $G_X \Rightarrow Pr(Y|do(X), Z, W) = Pr(Y|do(X), W)$.

- **Rule 2:** Action/observation exchange:

  $Y \perp\!\!\!\perp Z|(X, W)$ under $G_{\overline{X},\underline{Z}} \Rightarrow Pr(Y|do(X), do(Z), W) = Pr(Y|do(X), Z, W)$.

- **Rule 3:** Insertion/deletion of actions:

  $Y \perp\!\!\!\perp Z|(X, W)$ under $G_{\overline{X},\overline{Z(W)}} \Rightarrow Pr(Y|do(X), do(Z), W) = Pr(Y|do(X), W)$, where *Z(W)* is the set of *Z*-nodes that are not ancestors of any *W*-node in $G_{X^-}$.

Application of the do-calculus to a DAG entails several distinct steps:

- **Step 1:** Define an empirical model of interest expressed by a DAG *G*.

- **Step 2:** Generate a range of DAGs derived from the original DAG *G* according to the do-calculus graphical operations, e.g., $G_X^-$ or $G_{\underline{X}}$.

- **Step 3:** Generate the conditional independence relations associated with the DAGs generated from the subgraphs in step 2. These conditional independence relationships are obtained through the application of the Local Markov Condition and Graphoid Axioms to each DAG.

- **Step 4:** Check if one (or more) of the three do-calculus rules apply for a selected DAG and a selected conditional independence relationship.

- **Step 5:** If they apply, compute the do-operator equality associated with the rule in question. This step *defines* casual parameters.

- **Step 6:** *Identification* of a causal parameter occurs if the generated do-operator expressions for the causal parameters can be expressed in terms of the conditional distribution of observed variables.

These steps are summarized in Figure 1.

Goth (2006) and Huang and Valtorta (2006) show that Rules 1–3 are complete, meaning that they are *sufficient* for deriving all causal effects that can be identified by the conditional independence assumptions implicit in the DAG.

The do-calculus rules can be applied to any DAG, but they are restricted to models that can be defined as DAGs. In other words, the do-calculus rules only apply to models that can be

solely expressed by DAGs. They are not suitable for the identification of models whose assumptions lay outside the DAG conditional independence assumptions. In particular, they cannot identify the standard instrumental variable model examined in the next section.

We illustrate the use of the do-calculus through two examples. Section 4.1 examines the standard instrumental variable model. Section 4.2 examines the Front-Door Model.

## 4.1 The Instrumental Variable Model

We consider the simplest instrumental variable model that consists of four variables: (1) a confounding variable $U$ that is external and unobserved; (2) an external instrumental variable $Z$; (3) an observed variable $X$ caused by $U$ and $Z$; and (4) an outcome $Y$ caused by $U$ and $X$. The empirical instrumental variable model is described in the first column of Table 4. Its hypothetical counterpart is presented in the second column of Table 4.

The next step in the application of the do-calculus is the creation of DAGs derived from the original Instrumental Variable DAG. The rules of the do-calculus do not offer guidance on which DAGs to generate. For sake of brevity, we will consider four DAGs that can be operated on by the three do-calculus rules and thereby generate equalities involving the do-operator.

We first focus on the DAG $G_{\overline{X}}$ represented in the first DAG of Table 5. In this modified DAG, LMC (5) applied to variable $Z$ generates $Z \perp\!\!\!\perp (U, X, Y)$. Thereby $Y \perp\!\!\!\perp Z|X$ holds from application of the Symmetry and Weak Union properties of the Graphoid Axioms. We next map this relationship into Rule 1. Namely, if $Y \perp\!\!\!\perp Z|X$ holds in the Graph $G_{\overline{X}}$, then by Rule 1,

$$\Pr(Y|do(X), Z) = \Pr(Y|do(X)). \quad (11)$$

Equation (11) states that $Z$ is independent of $Y$ when $X$ is fixed. This implication of Rule 1 is captured by the relationship $Y \perp\!\!\!\perp Z|\tilde{X}$ in the hypothetical model. This relationship is automatically obtained by applying LMC (5) to variable $Z$ in the hypothetical model.

Using the Haavelmo approach, we create a hypothetical model shown in column 2 of Table 4. Without using any of Pearl's rules, it is immediate, applying LMC to $Z$, that

$$Y \perp\!\!\!\perp Z|\tilde{X}$$

and the causal parameter $Pr(Y|\tilde{X})$ is well defined. This parameter can be identified using standard IV methods widely used in econometrics (see, e.g., Matzkin, 2013).

We now use the do-calculus to show that causal effects can be obtained by conditioning on the unobserved confounding variable $U$. To illustrate the application of Rule 2, consider the graph $G_{\overline{X}, \underline{U}}$ depicted in DAG 2 of Table 5. In this modified DAG, LMC (5) applied to

variable $Y$ generates $Y \perp\!\!\!\perp (U, Z)|X$ and therefore $Y \perp\!\!\!\perp U|X$. If $Y \perp\!\!\!\perp U|X$ holds in the Graph $G_{\overline{X},\underline{U}}$ then by Rule 2,

$$\Pr(Y|do(X), do(U)) = \Pr(Y|do(X), U). \quad (12)$$

Next consider the graph $G_{\overline{U},\underline{X}}$ depicted in DAG 3 of Table 5. In this modified DAG, LMC (5) applied to variable $Y$ generates $Y \perp\!\!\!\perp (Z, X)|U$ and therefore $Y \perp\!\!\!\perp X|U$. Now if $Y \perp\!\!\!\perp X|U$ holds in the Graph $G_{\overline{U},\underline{X}}$ then by Rule 2,

$$\Pr(Y|do(X), do(U)) = \Pr(Y|X, do(U)). \quad (13)$$

Finally, consider the graph $G_{\underline{U}}$ depicted in DAG 4 of Table 5. Again, LMC (5) applied to variable $Y$ in this modified DAG generates $Y \perp\!\!\!\perp (U, Z)|X$ and therefore $Y \perp\!\!\!\perp U|X$. If $Y \perp\!\!\!\perp U|X$ holds in graph $G_{\underline{U}}$ by Rule 2,

$$\Pr(Y|X, do(U)) = \Pr(Y|X, U). \quad (14)$$

Combining Equations (12)–(14) we obtain:

$$\Pr(Y|do(X), U) = \Pr(Y|do(X), do(U)) = \Pr(Y|X, do(U)) = \Pr(Y|X, U). \quad (15)$$

Equation (15) states that the distribution of $Y$ conditioned on $U$ when $X$ is fixed is the same as the distribution of $Y$ conditioned on $U$ and $X$. In other words, the causal effects of $X$ on $Y$ can be obtained by conditioning $Y$ on the unobserved variable $U$. This well-known fact is easily obtained through the Hypothetical Model: LMC (5) on $X$ generates $Y \perp\!\!\!\perp X|(X, \tilde{U})$, and by Matching C-2, $Pr_H(Y|X, \tilde{U}) = Pr_E(Y|X, U)$.

As noted in Pearl (2009, chapters 3 and 5), the relationships discussed here — which exhaust the implications of the do-calculus for the instrumental variable model — are not sufficient to identify it unless the analyst can condition on $U$. Using the Haavelmo approach, a variety of identification strategies can identify the causal parameters, using instrumental variables. These strategies, however, rely on assumptions that cannot be expressed by DAGs. This renders the do-calculus unsuitable for the identification analysis of these type of models. Section 5 discusses this topic in detail. Heckman and Pinto (2013) use the hypothetical model approach to examine the necessary and sufficient conditions that identify the instrumental variable model.

### 4.2 Identifying the "Front-Door" Model

We next compare the do-calculus and an analysis based on our hypothetical model by identifying the causal effects of Pearl's "Front-Door model". Table 7 presents an example of the application of the rules of the do-calculus for the Front-Door model described in Table 6. That model consists of four variables: (1) an external unobserved variable $U$; (2) an observed variable $X$ caused by $U$; (3) an observed variable $M$ caused by $X$; and (4) an

outcome *Y* caused by *U* and *M*. The Front-Door model is presented in the first column of Table 6.

We are interested in identifying the distribution of the outcome *Y* when *X* is fixed at *x*. Within the context of the do-calculus, by identification we mean expressing the quantity *Pr(Y|do(X))* in terms of the distribution of observed variables.

The do-calculus identifies *Pr(Y|do(X))* through four steps, which we now perform. Steps 1, 2, and 3 identify *Pr(M|do(X))*, *Pr(Y|do(M))*, and *Pr(Y|M, do(X))*, respectively. Step 4 uses the first three steps to identify *Pr(Y|do(X))*.

1. Invoking LMC (5) for variable *M* of DAG $G_{\underline{X}}$, (DAG 1 of Table 7) generates $X \perp\!\!\!\perp M$. Thus, by Rule 2 of the do-calculus, we obtain *Pr(M|do(X)) = Pr(M|X)*.

2. Invoking LMC (5) for variable *M* of DAG $G_{\overline{M}}$, (DAG 1 of Table 7) generates $X \perp\!\!\!\perp M$. Thus, by Rule 3 of the do-calculus, *Pr(X|do(M)) = Pr(X)*. In addition, applying LMC (5) for variable *M* of DAG $G_{\underline{M}}$, (DAG 2 of Table 7) generates $M \perp\!\!\!\perp Y|X$. Thus, by Rule 2 of the do-calculus, *Pr(Y|X, do(M)) = Pr(Y|X, M)*.

$$
\begin{aligned}
\text{Therefore,} \Pr(Y|do(M)) &= \sum_{x' \in \mathrm{supp}(X)} \Pr(Y|X=x', do(M))\Pr(X=x'|do(M)) \\
&= \sum_{x' \in \mathrm{supp}(X)} \Pr(Y|X=x', M)\Pr(X=x'),
\end{aligned}
$$

where "supp" means support.

3. Invoking LMC (5) for variable *M* of DAG $G_{\overline{X},\underline{M}}$, (DAG 3 of Table 7) generates $Y \perp\!\!\!\perp M|X$. Thus, by Rule 2 of the do-calculus, *Pr(Y|M, do(X)) = Pr(Y|do(M), do(X))*. In addition, applying LMC (5) for variable *X* of DAG $G_{X,\overline{M}}$, (DAG 4 of Table 7) generates $(Y, M, U) \perp\!\!\!\perp X$. By weak union and decomposition, we obtain $Y \perp\!\!\!\perp X|M$. Thus, by Rule 3 of the do-calculus, we obtain that *Pr(Y|do(X), do(M)) = Pr(Y|do(M))*. Thus, *Pr(Y|M, do(X)) = Pr(Y|do(M), do(X)) = Pr(Y|do(M))*.

4. We collect the results from the three previous steps to identify *Pr(Y|do(X))* from observed data:

$$
\begin{aligned}
&\Pr(Y|do(X)=x) \\
&= \sum_{m \in \mathrm{supp}(M)} \Pr(Y|M, do(X)=x)\Pr(M|do(X)=x) \\
&= \sum_{m \in \mathrm{supp}(M)} \underbrace{\Pr(Y|do(M)=m, do(X)=x)}_{\text{Step 3}} \Pr(M=m|do(X)=x) \\
&= \sum_{m \in \mathrm{supp}(M)} \underbrace{\Pr(Y|do(M)=m)}_{\text{Step 3}} \Pr(M=m|do(X)=x) \\
&= \sum_{m \in \mathrm{supp}(M)} \underbrace{\left( \sum_{x' \in \mathrm{supp}(X)} \Pr(Y|X=x', M)\Pr(X=x') \right)}_{\text{Step 2}} \underbrace{\Pr(M=m|X=x)}_{\text{Step 1}}.
\end{aligned}
$$

In this fashion, we can use the do-calculus to identify the desired causal parameter. It is instructive to compare this proof of identification with one based on the approach of Haavelmo. We identify the causal effects of $X$ on $Y$ for the Front-Door model using a hypothetical model. We replace the relationship of $X$ on $M$ by a hypothetical variable $\tilde{X}$ that causes $M$. We use $Pr_E$ to denote the probability of the Front-Door model that generates the data (Column 1 of Table 6) and $Pr_H$ for the hypothetical model (Column 2 of Table 6). As before, we seek to identify $Pr_H(Y|\tilde{X})$ (the equivalent of $Pr(Y|do(X))$) from observed distributions in the empirical model.

We first present a lemma that states three useful conditional independence relations of the hypothetical model. The lemma is based on the application of LMC (5) and the Graphoid relationships:

**Lemma L-1.** In the Front-Door hypothetical model, (1) $Y \per\!\!\!\perp \tilde{X}|M$, (2) $\tilde{X} \per\!\!\!\perp M$, and (3) $Y \per\!\!\!\perp \tilde{X}|(M, X)$

*Proof.* By LMC (5) for $X$, we obtain $(Y, M, X) \per\!\!\!\perp \tilde{X}|U$. By LMC (5) for $Y$ we obtain $Y \per\!\!\!\perp (\tilde{X}, X)|(M, U)$. By Contraction applied to $(Y, M, X) \per\!\!\!\perp \tilde{X}|U$ and $Y \per\!\!\!\perp (\tilde{X}, X)|(M, U)$ we obtain $(Y, X) \per\!\!\!\perp \tilde{X}|(M, U)$. By LMC (5) for $U$ we obtain $(M, X) \per\!\!\!\perp \tilde{U}$. By Contraction applied to $(M, X) \per\!\!\!\perp U$ and $(Y, M, X) \per\!\!\!\perp \tilde{X}|U$ we obtain $(\tilde{X}, U) \per\!\!\!\perp (M, X)$. The second relationship in the Lemma is obtained by Decomposition. In addition, by Contraction on $(Y, X) \per\!\!\!\perp \tilde{X}|(M, U)$ and $(M, X) \per\!\!\!\perp \tilde{U}$ we obtain $(Y, X, U) \per\!\!\!\perp \tilde{X}|M$. The two remaining conditional independence relationships of the Lemma are obtained by Weak Union and Decomposition.[15]

Applying these results,

---

[15]One can also prove Lemma L-1 using Pearl's *d-Separation* criteria. According to Pearl (2009), a path $p$ connecting $X$ and $Y$ is said to be *d-Separated* (or blocked) by a set of nodes $Z$ if and only if

1.  a path $p$ contains a chain $i \rightarrow m \rightarrow j$ or a fork $i \leftarrow m \leftarrow j$ such that the middle node $m$ is in $Z$, or

2.  a path $p$ contains an inverted fork (or collider) $i \rightarrow m \leftarrow j$ such that the middle node $m$ is not in $Z$ and such that no descendant of $m$ is in $Z$.

A set $Z$ is said to *d-separate* $X$ from $Y$ if and only if $Z$ blocks every path from a node in $X$ to a node in $Y$. If $X$ and $Y$ are d-Separated by $Z$ according to a graph $G$, then $Y \perp\!\!\!\perp Y|Z$ in $G$. We are examining the Hypothetical Model described by second column of Table 2. Variables $Y$ and $X$ are connected by a single path $X \rightarrow M \rightarrow Y$. Thus we have that $Y \perp\!\!\!\perp \tilde{X}|M$, according to part 1 of the d-Separation criteria. Moreover, we can also state that $Y \perp\!\!\!\perp \tilde{X}|(M, X)$ as $X$ is not a collider nor a descendant of a collider (part 2 of the d-Separation criteria). Finally, there is no path that connects $X$ and $M$ of the form $X \rightarrow \dots \rightarrow M$ nor $X \leftarrow \dots \leftarrow M$. Thus we can state that $\tilde{X} \perp\!\!\!\perp M$ according to part 1 of the d-Separation criteria.

$$\mathrm{Pr}_{\mathrm{H}}(Y|\tilde{X}{=}x)$$

$$= \sum_{m\in\mathrm{supp}(M)} \mathrm{Pr}_{\mathrm{H}}(Y|M{=}m, \tilde{X}{=}x)\mathrm{Pr}_{\mathrm{H}}(M{=}m|\tilde{X}{=}x)$$

$$= \sum_{m\in\mathrm{supp}(M)} \mathrm{Pr}_{\mathrm{H}}(Y|M{=}m)\mathrm{Pr}_{\mathrm{H}}(M{=}m|\tilde{X}{=}x)$$

$$= \sum_{m\in\mathrm{supp}(M)} \left( \sum_{x'\in\mathrm{supp}(M)} \mathrm{Pr}_{\mathrm{H}}(Y|X{=}x', M{=}m)\mathrm{Pr}_{\mathrm{H}}(X{=}x'|M{=}m) \right) \mathrm{Pr}_{\mathrm{H}}(M{=}m|\tilde{X}{=}x)$$

$$= \sum_{m\in\mathrm{supp}(M)} \left( \sum_{x'\in\mathrm{supp}(X)} \mathrm{Pr}_{\mathrm{H}}(Y|X{=}x', M{=}m)\mathrm{Pr}_{\mathrm{H}}(X{=}x') \right) \mathrm{Pr}_{\mathrm{H}}(M{=}m|\tilde{X}{=}x)$$

$$= \sum_{m\in\mathrm{supp}(M)} \left( \sum_{x'\in\mathrm{supp}(X)} \mathrm{Pr}_{\mathrm{H}}(Y|X{=}x', \tilde{X}{=}x', M{=}m)\mathrm{Pr}_{\mathrm{H}}(X{=}x') \right) \mathrm{Pr}_{\mathrm{H}}(M{=}m|\tilde{X}{=}x)$$

$$= \sum_{m\in\mathrm{supp}(M)} \left( \sum_{x'\in\mathrm{supp}(X)} \underbrace{\mathrm{Pr}_{\mathrm{E}}(Y|M, X{=}x')}_{\text{by Theorem T-2}} \underbrace{\mathrm{Pr}_{\mathrm{E}}(X{=}x')}_{\text{by Theorem T-1}} \right) \underbrace{\mathrm{Pr}_{\mathrm{E}}(M{=}m|X{=}x)}_{\text{by Matching C-2}}.$$

The second equality comes from relationship (1) $Y \perp\!\!\!\perp X|\tilde{M}$ of Lemma L-1. The fourth equality comes from relationship (2) $X \perp\!\!\!\perp M$ of Lemma L-1. The fifth equality comes from relationship (3) $Y \perp\!\!\!\perp X|(\tilde{M}, X)$ of Lemma L-1. The last equality links the distributions of the hypothetical model with the ones of the empirical model. The first term uses Theorem T-2 to equate $Pr_{\mathrm{H}}(Y|X = x', \tilde{X} = x', M = m) = Pr_{\mathrm{E}}(Y|M, X = x')$. The second term uses the fact that $X$ is not a child of $\tilde{X}$; thus, by Theorem T-1, $Pr_{\mathrm{H}}(X = x') = Pr_{\mathrm{E}}(X = x')$. Finally, the last term uses Matching applied to $M$. Namely, LMC (5) for $M$ generates $M \perp\!\!\!\perp \tilde{X}|X$ in the hypothetical model. Then, by Matching C-2, $Pr_{\mathrm{H}}(M|\tilde{X} = x) = Pr_{\mathrm{E}}(M|X = x)$.

It is clear from the analysis presented in this section that, even when the do-calculus and the hypothetical model produce the same identification formulas, they differ greatly. Using the hypothetical model and adjoining $\tilde{X}$ to the original empirical model, we can define causal parameters in a straightforward way using the standard tools of probability theory. Analyses based on $\tilde{X}$ offer a simple, transparent, way to introduce Haavelmo's notion of fixing into models.

## 5 The Benefits and Limitations of DAGs

A major benefit of DAGs is their intuitively appealing description of models as causal chains. DAG assumptions list the variables in a model and their causal relationships. A DAG does not generate or characterize any restrictions on functional forms or parametric specifications of the structural equations. In this sense, if an identification result is achieved, it is obtained under very weak conditions.

This benefit of DAGs is also the source of their limitations. Methods that focus on identification of models solely described by DAGs lack the tools for invoking additional assumptions that would generate the identification of an *a priori* non-identified model. The instrumental variable model examined in Section 4.1 is a fundamental ingredient of a huge literature on econometric identification (see, e.g., Matzkin, 2013). It is the basis for identifying more sophisticated models such as the Generalized Roy model of Section 3.2,

which is widely used in econometrics in the analysis of selection bias and in evaluating social programs (Heckman, 1976, Heckman, 1979, Heckman and Robb, 1985, Powell, 1994, Heckman and Vytlacil, 2007a, b). Examples of this literature are nonparametric control functions (see, e.g., Blundell and Powell, 2003) and identification through instrumental variables (Reiersöl, 1945). Heckman and Pinto (2013) use the concept of a hypothetical model to examine the identification of the instrumental variable model. They present a unified approach that summarizes a range of identification strategies. Section 4.1 shows that the instrumental variable model is not identified applying the rules of the do-calculus. As mentioned, it is impossible to identify the causal effect of $X$ on $Y$ without using additional information.

The non-identification of the instrumental variable model poses a major limitation for the identification literature that relies exclusively on DAGs. Identification of the instrumental variable model relies on assumptions outside the scope of the DAG literature. For example, we can use LMC (5) to obtain the following conditional independence relationships: $Y \perp\!\!\!\perp Z|$ $(U, X)$ and $U \perp\!\!\!\perp Z$. These relationships in addition to $X \not\perp\!\!\!\perp Z$ satisfy the necessary criteria to apply the method of Two Stage Least Squares (TSLS). TSLS identifies the instrumental variable model under a linearity assumption. As a consequence, if we assume that the causal relationship of $X$ and $U$ on outcome $Y$ are represented by a linear equation, i.e., $Y = X\beta + U$, then it is well-known that parameters $\beta$ can be identified using $\text{cov}(Z, Y) / \text{cov}(Z, X)$ under standard rank conditions.

Linearity and homogeneity of the effects of $X$ on $Y$ across agents (i.e., $\beta$ is the same across the values $X, U$ take) are strong assumptions about the causal links that govern the relationship between $Y$ and $X$. This assessment fostered a huge literature in economics devoted to methods that relax linearity and homogeneity and that allow coefficients to be correlated with regressors. Examples of this literature are Imbens and Angrist (1994), Vytlacil (2002), and Heckman and Vytlacil (2005, 2007a, b), who identify the instrumental variable model under more general conditions by making assumptions on the causal relationship of $Z$ with $X$. Imbens and Angrist (1994) show that the instrumental variable model can be identified under a "monotonicity" assumption (increasing the values of an instrument has the same qualitative effect on all agents). Vytlacil (2002) shows that this assumption is equivalent to assuming an instrumental variable model in which the treatment assignment decision rule is separable in terms of unobserved characteristics of the agents and the instrumental variable. Heckman and Vytlacil (1999, 2005, 2007a, b) develop and apply this result.

Table 8 summarizes the common and distinct features of Pearl's do-calculus and the approach based on Haavelmo's hypothetical model. Both approaches use structural equation models in the sense of Koopmans and Reiersøl (1950). Both invoke autonomy and assume mutually independent errors $\varepsilon$. In recursive models, both use the Local Markov Condition and the Graphoid axioms. Both use "fixing" or the "do operator" to define counterfactuals.

The approaches diverge in their analyses of identification. The approach based on Haavelmo creates a hypothetical variable $\tilde{X}$ and an associated hypothetical model that is "outside the box" of the empirical model. It applies standard probability calculus to the hypothetical

model to connect the hypothetical model to the empirical model. Pearl's do-calculus creates a new set of extra-statistical tools to identify the causal parameters created by fixing or the "do-operator." Our analysis shows that in the hypothetical model of Haavelmo, the special extra-statistical tools of the do-calculus are not required to identify causal parameters. By relying exclusively on statistical tools, the hypothetical model approach allows for additional econometric assumptions that identify a broader range of models that cannot be identified using the rules that only apply to DAGs, i.e. the "do-calculus."

## 6 Hypothetical Models and Simultaneous Equations

The literature on causality provides a framework for modeling causal processes that are based on DAGs. Less is known about Directed Cyclic Graphs (DCGs) that are used to represent Simultaneous Equations. Indeed, the fundamental Local Markov Condition no longer holds for DCGs (Spirtes, 1995). Nevertheless, the notion of fixing readily extends to a system of simultaneous equations.

Consider a system of two equations:

$$Y_1 = g_{Y_1}(Y_2, X_1, U_1), \quad \text{(7a)}$$

$$Y_2 = g_{Y_2}(Y_1, X_2, U_2). \quad \text{(7b)}$$

$\mathcal{T}_E = \{Y_1, Y_2, X_1, X_2, U_1, U_2\}$. Our analysis can be readily generalized to systems with more than two equations, but for the sake of brevity, we focus on the two-equation case. To simplify notation, we keep the variables in $\varepsilon$ implicit.

The empirical Simultaneous Equations Model of (7a) and (7b) is represented as Model 1 of Table 9. Many different versions of this model appear in the literature. For simplicity, we assume $U_1 \perp\!\!\!\perp U_2$ and $(U_1, U_2) \perp\!\!\!\perp (X_1, X_2)$.[16]

The hypothetical model associated with the causal operation of fixing both $Y_2$ and $Y_1$ is represented in Model 2 of Table 9. Under autonomy, the causal effect of $Y_2$ on $Y_1$ when $Y_2$ is fixed at $y_2$ is given by $Y_1(y_2) = g_{Y_1}(y_2, X, U_1)$. Symmetrically, $Y_2(y_1) = g_{y2}(y_1, X, U_2)$. We define hypothetical random variables $\tilde{Y}_1, \tilde{Y}_2$. They replace the $Y_1, Y_2$ inputs in Equations (7a) and (7b) in the same fashion as discussed in previous sections. $(\tilde{Y}_1, \tilde{Y}_2) \perp\!\!\!\perp (X_1, X_2, U_1, U_2)$; and $\tilde{Y}_1 \perp\!\!\!\perp \tilde{Y}_2$. $\mathcal{T}_H = \{\tilde{Y}_1, \tilde{Y}_2, \tilde{Y}_1, \tilde{Y}_2, X_1, X_2, U_1, U_2\}$. We assume a common support for $(Y_1, Y_2)$ and $(\tilde{Y}_1, \tilde{Y}_2)$.

In the same fashion as in the model previously discussed, the distribution of $Y_1$ when $Y_2$ is fixed at $y_2$ is given by $Pr_H(Y_1|Y_2 = y_2)$. The average causal effect of $Y_2$ on $Y_1$ when $Y_2$ is fixed at the two values of $y_2$ and $y'_2$ is given by $\mathbf{E}_H(Y_1|\tilde{Y}_2 = y_2) - \mathbf{E}_H(Y_1|\tilde{Y}_2 = y'_2)$, where $\mathbf{E}_H$ denotes expectation over the probability measure $Pr_H$ of the hypothetical model. The

---

[16]These assumptions are made to simplify the analysis. A large literature relaxes these assumptions and develops identification criteria for cases where $U_1 \not\perp\!\!\!\perp U_2$ and $(U_1, U_2) \not\perp\!\!\!\perp (X_1, X_2)$. The literature considers a variety of specifications (see Matzkin, 2008). We maintain the assumptions that $U_1 \perp\!\!\!\perp U_2$ and $(U_1, U_2) \perp\!\!\!\perp (X_1, X_2)$ for simplicity.

hypothetical variation of $\tilde{Y}_2$ corresponds to the standard Marshallian and Walrasian thought experiments in which quantities or prices are fixed to trace out demand and supply curves (see, e.g., Mas-Colell et al., 1995). A symmetric analysis produces the causal effect of $Y_1$ on $Y_2$. Thus we obtain the counterpart to the counterfactuals defined for the recursive models earlier in this paper.

Under simultaneity, the graph for Model 1 is cyclic and the relationships that hold for DAGs, such as the LMC (5), break down (Lauritzen and Richardson, 2002; Spirtes, 1995). Equations (7a) and (7b) cannot be represented as Directed Bayesian networks. The tools developed for DAGs do not directly apply and require modification. Equations (7a) and (7b) are fundamentally non-recursive and observed variables emerge from a feedback process.

The do-calculus does not have the concept of a hypothetical model. It focuses exclusively on the empirical model. Thus, in the simultaneous equation models depicted in Table 9, the do-calculus can only be applied to the model in the first panel (the empirical model). In the case of the simultaneous equation model, the DAG is cyclical and the rules of the do-calculus do not hold.

A traditional assumption in the simultaneous equations literature is "completeness"—the existence of at least a local solution for $Y_1$ and $Y_2$ in terms of $(X_1, X_2, U_1, U_2)$:

$$Y_1 = \phi_1(X_1, X_2, U_1, U_2), \quad \text{(8a)}$$

$$Y_2 = \phi_2(X_1, X_2, U_1, U_2).\text{[17]} \quad \text{(8b)}$$

These are called "reduced form" equations (see, e.g., Matzkin, 2008, 2013). They inherit the autonomy properties of the structural equations.[18]

The assumption of the existence of a reduced form is not innocuous even in the linear cases for continuous $Y_1$ and $Y_2$ analyzed by Haavelmo (1943, 1944) and the Cowles Foundation pioneers (see Koopmans et al., 1950). Heckman (1978), Tamer (2003), and Chesher and Rosen (2012) analyze the case in which $Y_1$ and $Y_2$ are discrete valued. Solutions (8a) and (8b) may not exist except under conditions given in those papers.[19] Alternatively, there may be multiple solutions giving rise to reduced form correspondences. In the case where no solutions exist, the model is incoherent as an equilibrium model unless additional assumptions are invoked. However, one can construct hypothetical models using Haavelmo's insights even in incoherent cases.[20]

In addition, some frameworks for multivariate discrete data may not be sufficiently rich to distinguish correlation from causation. Heckman (1978) shows that log-linear models for discrete data used in statistics (see, e.g., Bishop et al., 1975) have too few parameters to make causal distinctions. He introduces a class of latent variable models in which such distinctions are possible.

Note further that even in models in which the reduced form equations are well defined, it is not possible, in general, to *simultaneously* vary $\tilde{Y}_1$ and $\tilde{Y}_2$ so that they (i) solve Equations

(7a) and (7b) and (ii) also satisfy the requirement that ($\tilde{Y}_1$, $\tilde{Y}_2$) $\perp\!\!\!\perp$ ($X_1$, $X_2$, $U_1$, $U_2$). This is apparent from the reduced form equations (8a) and (8b) that, under completeness, the proposed variations must also satisfy. Nonetheless, $\tilde{Y}_2$ and $\tilde{Y}_1$ can be separately constructed to create hypothetical models corresponding to Equations (7a) and (7b) respectively. These equations exist as theoretical constructs independent of any particular equilibrium construct.

Matzkin (2007, 2008, 2012, 2013) presents comprehensive and definitive treatments of alternative approaches for identifying simultaneous equations. Our analysis readily extends to systems with more than two equations, but for the sake of brevity we do not make the extension here.

## 7 Summary and Conclusions

This paper examines Haavelmo's fundamental contributions to the study of causal inference. He produced the first formal analysis of the distinction between causation and correlation. He carefully distinguished the process of defining causality—a mental act that assigns hypothetical variation to inputs—from the act of identifying causal models from data. Haavelmo was remarkably clear about concepts that are still muddled in some quarters of statistics.[21]

Haavelmo shows us that causal effects of inputs on outputs are defined in abstract models that assign independent variation to inputs. He formalized Frisch's notion that causality is in

---

[18]Under completeness, we can use a version of indirect least squares to define causal parameters and identify them where the induced variation in $\tilde{Y}_1$ and $\tilde{Y}_2$ satisfy equilibrium conditions. Thus if $X_1$ and $X_2$ are disjoint, one can use indirect least squares to identify from reduced form equations (8a) and (8b), assumed to be differentiable:

$$\underset{\text{(From 8a)}}{\frac{\partial Y_1}{\partial X_2}} = \frac{\partial g_{Y_1}(Y_2, X_1, U_1)}{\partial Y_2} \quad \underset{\text{(From 8b)}}{\frac{\partial Y_2}{\partial X_2}} \quad .$$

From the reduced forms:

$$\frac{\partial Y_1}{\partial X_2} = \frac{\partial \phi_1(\cdot)}{\partial X_2},$$
$$\frac{\partial Y_2}{\partial X_2} = \frac{\partial \phi_2(\cdot)}{\partial X_2}.$$

Thus,

$$\frac{\frac{\partial Y_1}{\partial X_2}}{\frac{\partial Y_2}{\partial X_2}} = \frac{\frac{\partial \phi_1(\cdot)}{\partial X_2}}{\frac{\partial \phi_2(\cdot)}{\partial X_2}} = \frac{\partial g_{Y_1}(Y_2, X_1, U_1)}{\partial Y_2} \quad .$$

If $X_1$ and $X_2$ contain common elements, the method can be modified to use only the distinct elements in $X_1$ and $X_2$ in this analysis.
[19]Linear probability model approximations to Equations (7a) and (7b), as advocated by Angrist and Pischke (2008), although widely used, are in general not autonomous. They can, however, be estimated and identified for incoherent models, creating the illusion of coherency through approximation error. See Heckman and MaCurdy (1985) for a discussion.
[20]This might be a conceptually unsatisfactory exercise unless the data intended to be described by the model display disequilibrium cycling phenomena and a time sequence for the evolution of the system, e.g., $Y_1^{(t)}$, $Y_2^{(t+1)}$, ..., is postulated as functions of inputs where superscripts denote time-dated variables.
[21]See, e.g., Holland (1986) and Sobel (2005) for examples of the confusion between models and identification strategies exemplified by the claim that no causal statements are possible unless persons are randomly assigned to treatment.

the mind. We formalize his insight extending his analysis for linear models to more general models. This enables us to discuss causal concepts such as "fixing" using an intuitive approach that applies Haavelmo's ideas.

Following Haavelmo, we distinguish the definition of causal parameters from their identification. Our approach to defining causality relies on the assumption of autonomy joined with Haavelmo's notion of hypothetical random variables. Together they enable us to express the distribution of counterfactual outcomes using structural equations and the distributions of the data by replacing the variables whose causal effects we seek to establish with their hypothetical counterparts. Causal models thus defined apply standard statistical tools and do not require new procedures like the do-calculus that lie outside the scope of the standard tools of probability and statistics.

Identification in Haavelmo's model is achieved in recursive models by applying standard statistical tools to Bayesian Networks. We link the distributions of empirical and hypothetical models by expressing the quantities of interest in the hypothetical model into observed quantities in the empirical one.

We illustrate the benefits of Haavelmo's approach first by analyzing an IV DAG and then by comparing identification of the causal effects of Pearl's flagship Front-Door model (Pearl, 2009) using a method based on the Haavelmo approach and a method based on the do-calculus of Pearl (2009). The IV model is not identified using solely the calculus of the do-operator. While both methods generate the same estimator for the Front-Door model, the identification methods differ on both conceptual and methodological grounds. The definition of causal parameters is clearly separated from their identification in the Haavelmo approach. The two tasks are often confounded in applications of the do-calculus. The do-calculus is cumbersome compared to the elegance and simplicity of the Haavelmo approach. We discuss the limitations of methods of identification that rely on the fundamentally recursive approach of Directed Acyclic Graphs.

That framework cannot accommodate the fundamentally non-recursive simultaneous equations model. The hypothetical model readily accommodates an analysis of causality in the simultaneous equations model. The framework of simultaneous equations is fundamentally non-recursive and falls outside of the framework of Bayesian causal nets and DAGs. The rigorous definition of causality in a variety of models including the simultaneous equations framework and the identification of causal parameters, are central and enduring contributions of Haavelmo (1944).

## Acknowledgments

## A Mathematical Appendix

Theorem T-1:

*Proof.* If $T$ is non-descendant of $\tilde{X}$ in the hypothetical model, i.e., $T \in \mathscr{T}_E \setminus D_H(\tilde{X})$, then $T \in \mathscr{T}_E \setminus Ch_H(\tilde{X})$ as $Ch_H(\tilde{X}) \subset D_H(\tilde{X})$. Thus, $Pr_H(T|Pa_H(T)) = Pr_E(T|Pa_E(T))$ from Equation (8). Moreover, it must be the case that parents of $T$ are also non-descendants of $\tilde{X}$; i.e., $Pa_H(T) \subset \mathscr{T}_E \setminus D_H(\tilde{X}) \subset \mathscr{T}_E \setminus Ch_H(\tilde{X}) \therefore Pa_H(T) = Pa_E(T)$ by Equation (8). Another way of saying this is that the parents of $T$ are not children of $\tilde{X}$. Thus, we can use factorization (6) to write:

$$\mathrm{Pr}_H(\mathscr{T}_E \setminus D_H(\tilde{X})) = \prod_{T \in \mathscr{T}_E \setminus D_H(\tilde{X})} \mathrm{Pr}_H(T|Pa_H(T)) = \prod_{T \in \mathscr{T}_E \setminus D_H(\tilde{X})} \mathrm{Pr}_E(T|Pa_E(T)) = \mathrm{Pr}_E(\mathscr{T}_E \setminus D_H(\tilde{X})).$$

As a consequence, $Pr_H(W) = Pr_E(W)$ for all $W \subset \mathscr{T}_E \setminus D_H(\tilde{X})$, and thereby

$$\mathrm{Pr}_H(W=w|Z=z) = \frac{\mathrm{Pr}_H(W=w, Z=z)}{\mathrm{Pr}_H(Z=z)} = \frac{\mathrm{Pr}_E(W=w, Z=z)}{\mathrm{Pr}_E(Z=z)} = \mathrm{Pr}_E(W=w|Z=z).$$

Conditioning on $\tilde{X}$ comes from that fact that $\tilde{X} \perp\!\!\!\perp (\mathscr{T}_E \setminus D_H(\tilde{X}))$, which is obtained by applying LMC (5) to $\tilde{X}$ in $G_H$.

Theorem T-2:

*Proof.* In order to prove the theorem, we first partition the set of variables $\mathscr{T}_E$ into four sets:

$$\mathscr{T}_E = \underbrace{\{\mathscr{T}_E \setminus D_E(X)\}}_{\text{Set 1}} \cup \underbrace{\{D_E(X) \setminus Ch_E(X)\}}_{\text{Set 2}} \cup \underbrace{\{Ch_H(X)\}}_{\text{Set 3}} \cup \underbrace{\{Ch_H(\tilde{X})\}}_{\text{Set 4}}.$$

Set 1 consists of all variables in $\mathscr{T}_E$ that are non-descendants of $X$ in the empirical model and thereby nondescendants of $\tilde{X}$ in the hypothetical one. Set 2 consists of descendants of $X$ but not directly caused by $X$, i.e., except its Children. Sets 3 and 4 are the Children of $X$ and $\tilde{X}$ in the hypothetical model. Note that Sets 3 and 4 consist of all Children of $X$ in the empirical model as $Ch_E(X) = Ch_H(X) \cup Ch_H(\tilde{X})$. We now examine the variables of each set separately:

1. For all $T \in \mathscr{T}_H \setminus D_E(X) \Rightarrow \{T, Pa_H(T)\} \subset \mathscr{T}_H \setminus D_E(X) \subset \mathscr{T}_E \setminus D_H(\tilde{X})$, as $D_H(\tilde{X}) \subset D_E(X)$. Also $X \in \mathscr{T}_E \setminus D_H(\tilde{X})$. Thus, by Theorem T-1, $Pr_H(T|Pa_H, \tilde{X} = x, X = x) = Pr_E(T|Pa_E(T), X = x)$.

2. $T \in D_E(X) \setminus Ch_E(X) \Rightarrow X \notin Pa_H(T), \tilde{X} \notin Pa_H(T)$, and $Pa_H(T) = Pa_E(T)$. Moreover, $X, \tilde{X}$ must be non-descendants of $T$ due to the acyclic property of the empirical model on $X$. Thus, by LMC (5), $(X, \tilde{X}) \perp\!\!\!\perp T|Pa_H(T)$. By Weak Union, $X \perp\!\!\!\perp T| (Pa_H(T), \tilde{X})$. Therefore $Pr_H(T|Pa_H(T), \tilde{X} = x, X = x) = Pr_H(T|Pa_H(T), X = x) = Pr_E(T|Pa_E(T), X = x)$ by Equation (8).

3. $T \in Ch_H(X) \Rightarrow X \notin Pa_H(T)$ and $X \in Pa_H(T) = Pa_E(T)$. Also, $\tilde{X}$ is external, thus $\tilde{X} \perp\!\!\!\perp T|Pa_H(T)$ by LMC (5) applied to $T$. Therefore $Pr_H(T|Pa_H(T) \setminus X, \tilde{X} = x, X = x) = Pr_H(T|Pa_H(T) \setminus X, X = x) = Pr_E(T|Pa_E(T) \setminus X, X = x)$ by Equation (8) as $T \in Ch_H(X) \subset \mathscr{T}_E \setminus Ch_H(X)$.

4. $T \in Ch_H(X) \Rightarrow \tilde{X} \in Pa_H(T)$. Moreover, $X$ must be a non-descendant of $T$ due to the acyclic property of the empirical model on $X$. Thus, by LMC (5), $X \perp\!\!\!\perp T|Pa_H(T)$. Therefore $Pr_H(T|Pa_H(T) \setminus X, \tilde{X} = x, X = x) = Pr_H(T|Pa_H(T) \setminus X, \tilde{X} = x) = Pr_E(T|Pa_E(T) \setminus X, X = x)$ by Equation (9).

Grouping items 1–4, we have that for all $T \in \mathscr{T}_H$, $Pr_H(T|Pa_H(T), \tilde{X} = x, X = x) = Pr_E(T|Pa_E(T), X = x)$. Thus we can use the factorization (6) to obtain

$$
\begin{aligned}
\mathrm{Pr}_H(\mathscr{T}_E|X=x, \tilde{X}=x) &= \prod_{T \in \mathscr{T}_E} \mathrm{Pr}_H(T|Pa_H(T), \tilde{X}=x, X=x) \\
&= \prod_{T \in \mathscr{T}_E} \mathrm{Pr}_E(T|Pa_E(T), X=x) \qquad (17) \\
&= \mathrm{Pr}_E(\mathscr{T}_E|X=x).
\end{aligned}
$$

The claim of the theorem is a direct consequence of Equation (17).

Corollary C-1:

*Proof.*

$$
\begin{aligned}
\mathrm{Pr}_H(\mathscr{T}_E|X=\tilde{X}) &= \sum_{x \in \mathrm{supp}(X)} \mathrm{Pr}_H(\mathscr{T}_E|X=x, \tilde{X}=x) \frac{\mathrm{Pr}_H(X=x,\tilde{X}=x)}{\sum_{x \in \mathrm{supp}(X)} \mathrm{Pr}_H(X=x,\tilde{X}=x)} \\
&= \sum_{x \in \mathrm{supp}(X)} \mathrm{Pr}_H(\mathscr{T}_E|X=x, \tilde{X}=x) \frac{\mathrm{Pr}_H(X=x)\mathrm{Pr}_H(\tilde{X}=x)}{\sum_{x \in \mathrm{supp}(X)} \mathrm{Pr}_H(X=x)\mathrm{Pr}_H(\tilde{X}=x)} \\
&= \sum_{x \in \mathrm{supp}(X)} \mathrm{Pr}_H(\mathscr{T}_E|X=x, \tilde{X}=x) \frac{\mathrm{Pr}_H(X=x)}{\sum_{x \in \mathrm{supp}(X)} \mathrm{Pr}_H(X=x)} \\
&= \sum_{x \in \mathrm{supp}(X)} \mathrm{Pr}_H(\mathscr{T}_E|X=x, \tilde{X}=x)\mathrm{Pr}_H \mathrm{Pr}_H(X=x) \\
&= \sum_{x \in \mathrm{supp}(X)} \mathrm{Pr}_E(\mathscr{T}_E|X=x)\mathrm{Pr}_E(X=x) \\
&= \mathrm{Pr}_E(\mathscr{T}_E).
\end{aligned}
$$

The second equality stems from $Pa_H(\tilde{X}) = \varnothing$ and $X$ is not descendant of $\tilde{X}$, thus by LMC (5), $X \perp\!\!\!\perp \tilde{X}$. Therefore $Pr_H(X = x, \tilde{X} = x) = Pr_H(X = x)Pr_H(\tilde{X} = x)$. The third equality comes from the assumption that $Pr_H(\tilde{X} = x)$ is constant due to uniformity. The fourth equality comes from the fact that $\sum_{x \in \mathrm{supp}(X)} \mathrm{Pr}_H(X=x)=1$. The first term of the fifth equality comes from an application of Theorem T-2. The second term of the fifth equality comes from Theorem T-1 and the fact that $X \in \mathscr{T}_E \setminus D_H(\tilde{X})$.

Theorem T-3:

*Proof.*

$$
\begin{aligned}
\mathrm{Pr}_{\mathrm{H}}(\mathscr{T}_{\mathrm{E}} \backslash X | \tilde{X}{=}x) &= \prod_{T \in \mathscr{T}_{\mathrm{E}} \backslash \{X \cup Ch_{\mathrm{H}}(X)\}} \mathrm{Pr}_{\mathrm{H}}(T|Pa(T)) \prod_{T \in Ch_{\mathrm{H}}(X)} \mathrm{Pr}_{\mathrm{H}}(T|Pa(T)\backslash \tilde{X}, \tilde{X}{=}x) \\
&= \prod_{T \in \mathscr{T}_{\mathrm{E}} \backslash \{X \cup Ch_{\mathrm{E}}(X)\}} \mathrm{Pr}_{\mathrm{H}}(T|Pa(T)) \prod_{T \in Ch_{\mathrm{E}}(X)} \mathrm{Pr}_{\mathrm{H}}(T|Pa(T)\backslash \tilde{X}, \tilde{X}{=}x) \\
&= \prod_{T \in \mathscr{T}_{\mathrm{E}} \backslash \{X \cup Ch_{\mathrm{E}}(X)\}} \mathrm{Pr}_{\mathrm{E}}(T|Pa(T)) \prod_{T \in Ch_{\mathrm{E}}(X)} \mathrm{Pr}_{\mathrm{E}}(T|Pa(T)\backslash X, X{=}x) \\
&= \mathrm{Pr}_{\mathrm{E}}(\mathscr{T}_{\mathrm{E}} \backslash X | do(X){=}x).
\end{aligned}
$$

The first equality comes from the fact that the hypothetical model is a DAG, therefore we apply factorization (6). The second equality comes from the characteristic of the do-operator, which targets all causal links of a fixed variable $X$. Thus, the hypothetical variable $\tilde{X}$ must replace all $X$ inputs which is equivalent to $Ch_{\mathrm{H}}(X) = Ch_{\mathrm{E}}(X)$. The first and second terms of the third equality come as a consequence of Equations (8) and (9), respectively. The last equality comes from the definition of the do-operator.

Matching C-2:

*Proof.*

$$
\begin{aligned}
\mathrm{Pr}_{\mathrm{H}}(W|Z, \tilde{X}{=}x) &= \mathrm{Pr}_{\mathrm{H}}(W|Z, \tilde{X}{=}x, X{=}x) \text{ by assumption } X \perp\!\!\!\perp W|(Z, \tilde{X}) \text{ in } G_{\mathrm{H}} \\
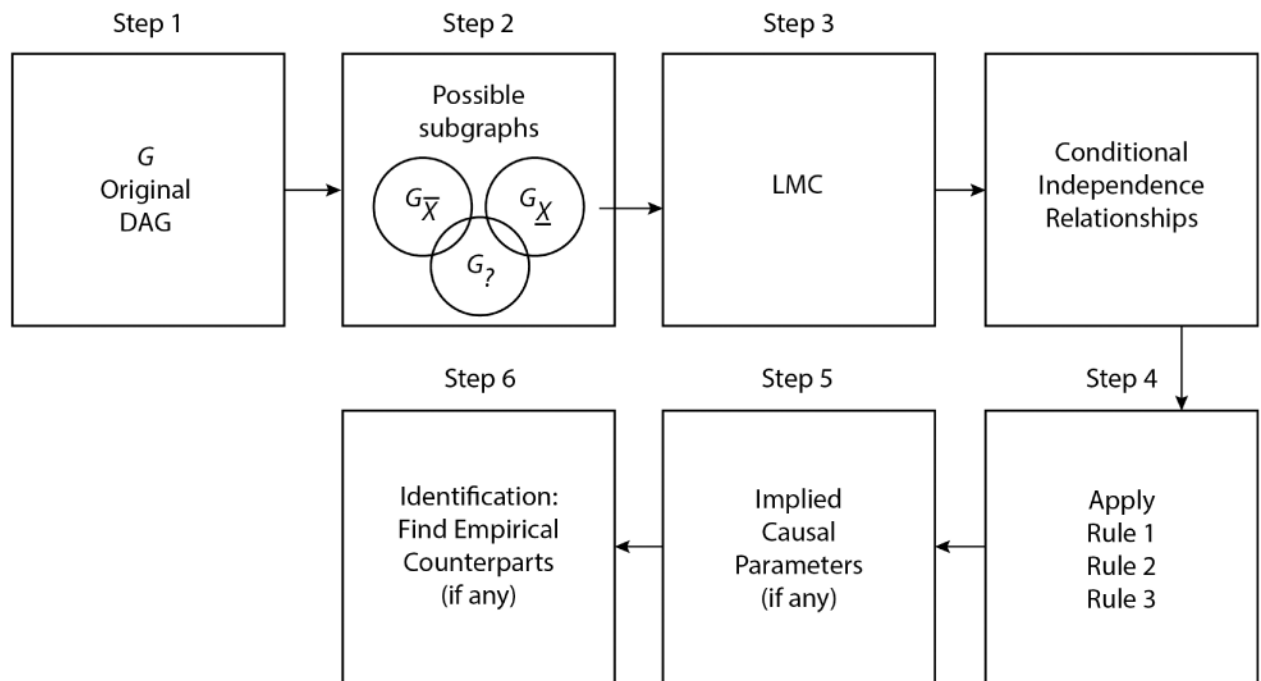&= \mathrm{Pr}_{\mathrm{E}}(W|Z, X{=}x) \text{ by Theorem } T-2.
\end{aligned}
$$

# References

Angrist, JD.; Pischke, JS. Mostly Harmless Econometrics: An Empiricist's Companion. Princeton: Princeton University Press; 2008.

Bishop, YM.; Fienberg, SE.; Holland, PW. Discrete Multivariate Analysis: Theory and Practice. Cambridge, Massachusetts: The MIT Press; 1975.

Blundell, R.; Powell, J. Endogeneity in nonparametric and semiparametric regression models. In: Dewatripont, LPHM.; Turnovsky, SJ., editors. Advances in Economics and Econometrics: Theory and Applications, Eighth World Congress. Vol. 2. Cambridge, UK: Cambridge University Press; 2003.

Chalak K, White H. Causality, conditional independence, and graphical separation in settable systems. Neural Computation. 2012; 24(7):1611–1668.

Chesher A, Rosen A. Simultaneous equations for discrete outcomes: Coherence, completeness, and identification. Working Papers CWP21/12, cemmap. 2012

Dawid A. Separoids: A mathematical framework for conditional independence and irrelevance. Annals of Mathematics and Artificial Intelligence. 2001; 32(1–4):335–372.

Dawid AP. Conditional independence in statistical theory (with discussion). Journal of the Royal Statistical Society Series B (Statistical Methodological). 1979; 41(1):1–31.

Fechner GT. Outline of a new principle of mathematical psychology. Psychological Research. 1851; 49:203–207. [PubMed: 3327074]

Freedman, D. Statistical Models and Causal Inference: A Dialogue with the Social Sciences. Collier, D.; Sekhon, J.; Stark, P., editors. Cambridge, UK: Cambridge University Press; 2010.

Frisch, R. A Dynamic Approach to Economic Theory: The Yale Lectures of Ragnar Frisch, 1930. Bjerkholt, O.; Qin, D., editors. New York, New York: Routledge; 1930. published 2010

Frisch, R. The Foundations of Econometric Analysis. Cambridge University Press; 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)

Goth G. Judea pearl interview. IEEE Internet Computing. 2006; 10(5):6.

Haavelmo T. The statistical implications of a system of simultaneous equations. Econometrica. 1943 Jan; 11(1):1–12.

Haavelmo T. The probability approach in econometrics. Econometrica. 1944; 12(Supplement):iii–vi. 1–115.

Hansen LP, Sargent TJ. Formulating and estimating dynamic linear rational expectations models. Journal of Economic Dynamics and Control. 1980 Feb; 2(1):7–46.

Heckman, J.; Pinto, R. A unified approach to examine treatment effects: Causality and identification. University of Chicago, Department of Economics; 2013. Unpublished manuscript

Heckman JJ. The common structure of statistical models of truncation, sample selection and limited dependent variables and a simple estimator for such models. Annals of Economic and Social Measurement. 1976 Dec; 5(4):475–492.

Heckman JJ. Dummy endogenous variables in a simultaneous equation system. Econometrica. 1978 Jul; 46(4):931–959.

Heckman JJ. Sample selection bias as a specification error. Econometrica. 1979 Jan; 47(1):153–162.

Heckman JJ. The scientific model of causality. Sociological Methodology. 2005 Aug; 35(1):1–97.

Heckman JJ. Econometric causality. International Statistical Review. 2008 Apr; 76(1):1–27.

Heckman JJ, MaCurdy TE. A simultaneous equations linear probability model. Canadian Journal of Economics. 1985 Feb; 18(1):28–37.

Heckman JJ, Robb R. Alternative methods for evaluating the impact of interventions: An overview. Journal of Econometrics. 1985 Oct-Nov;30(1–2):239–267.

Heckman JJ, Vytlacil EJ. Local instrumental variables and latent variable models for identifying and bounding treatment effects. Proceedings of the National Academy of Sciences. 1999 Apr; 96(8): 4730–4734.

Heckman JJ, Vytlacil EJ. Structural equations, treatment effects and econometric policy evaluation. Econometrica. 2005 May; 73(3):669–738.

Heckman, JJ.; Vytlacil, EJ. Econometric evaluation of social programs, part I: Causal models, structural models and econometric policy evaluation. In: Heckman, J.; Leamer, E., editors. Handbook of Econometrics. Vol. 6B. Amsterdam: Elsevier; 2007a. p. 4779-4874.

Heckman, JJ.; Vytlacil, EJ. 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: Heckman, J.; Leamer, E., editors. Handbook of Econometrics. Vol. 6B. Amsterdam: Elsevier; 2007b. p. 4875-5143.Chapter 71

Heidelberger, M. Nature from within: Gustav Theodor Fechner and his psychophysical worldview. Pittsburgh, PA: University of Pittsburgh Press; 2004.

Holland PW. Statistics and causal inference. Journal of the American Statistical Association. 1986 Dec; 81(396):945–960.

Howard, RA.; Matheson, JE. Influence diagrams. 1. Menlo Park, CA: Stanford Research Institute; 1981. Principles and applications of decision analysis; p. 720-762.

Huang, Y.; Valtorta, M. Technical report. University of South Carolina Department of Computer Science; 2006. A study of identifiability in causal Bayesian network.

Huang Y, Valtorta M. On the completeness of an identifiability algorithm for semi-markovian models. Annals of Mathematics and Artificial Intelligence. 2008; 54(4):363–408.

Imbens GW, Angrist JD. Identification and estimation of local average treatment effects. Econometrica. 1994 Mar; 62(2):467–475.

Kiiveri H, Speed TP, Carlin JB. Recursive causal models. Journal of the Australian Mathematical Society (Series A). 1984; 36(1):30–52.

Koopmans TC, Reiersøl O. The identification of structural characteristics. The Annals of Mathematical Statistics. 1950 Jun; XXI(2):165–181.

Koopmans, TC.; Rubin, H.; Leipnik, RB. Measuring the equation systems of dynamic economics. In: Koopmans, TC., editor. Statistical Inference in Dynamic Economic Models, Number 10 in Cowles Commission Monograph. Vol. Chapter 2. New York: John Wiley & Sons; 1950. p. 53-237.

Lauritzen, SL. Graphical Models. Oxford, UK: Clarendon Press; 1996.

Lauritzen, SL. Causal inference from graphical models. In: Barndorff-Nielsen, O.; Cox, DR.; Klüppelberg, C., editors. Complex Stochastic Systems. London: Chapman and Hall; 2001. p. 63-107.

Lauritzen SL, Richardson TS. Chain graph models and their causal interpretations. Journal of the Royal Statistical Society: Series B (Statistical Methodology). 2002; 64(3):321–348.

Lehmann, EL.; Romano, JP. Testing Statistical Hypotheses. Third. New York: Springer Science and Business Media; 2005.

Margolis, M.; List, J.; Osgood, D. Endangered options and endangered species: what we can learn from a dubious design. Gettysburg College, Department of Economics; 2012 Apr. Unpublished manuscript

Marshall, A. Principles of Economics. New York: Macmillan and Company; 1890.

Mas-Colell, A.; Whinston, MD.; Green, JR. Microeconomic Theory. New York: Oxford University Press; 1995.

Matzkin, RL. Nonparametric identification. In: Heckman, J.; Leamer, E., editors. Handbook of Econometrics. Vol. 6B. Amsterdam: Elsevier; 2007.

Matzkin RL. Identification in nonparametric simultaneous equations models. Econometrica. 2008; 76(5):945–978.

Matzkin RL. Identification in nonparametric limited dependent variable models with simultaneity and unobserved heterogeneity. Journal of Econometrics. 2012; 166(1):106–115.

Matzkin RL. Nonparametric identification of structural economic models. Annual Review of Economics. 2013; 5 Forthcoming.

Newey WK, Powell JL. Instrumental variable estimation of nonparametric models. Econometrica. 2003 Sep; 71(5):1565–1578.

Pearl, J. Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference. San Mateo, CA: Morgan Kaufmann Publishers Inc.; 1988.

Pearl J. [Bayesian analysis in expert systems]: Comment: Graphical models, causality and intervention. Statistical Science. 1993; 8(3):266–269.

Pearl J. Causal diagrams for empirical research. Biometrika. 1995 Dec; 82(4):669–688.

Pearl, J. Causality. Cambridge, England: Cambridge University Press; 2000.

Pearl, J. Causality: Models, reasoning, and inference. New York: Cambridge University Press; 2001. Reprinted with corrections ed.

Pearl, J. Causality: Models, Reasoning, and Inference. 2nd. New York: Cambridge University Press; 2009.

Pearl, J.; Verma, TS. A theory of inferred causation. In: Prawitz, D.; Skyrms, B.; Westerståhl, D., editors. Logic, Methodology, and Philosophy of Science; Proceedings of the Ninth International Congress of Logic, Methodology, and Philosophy of Science; Uppsala, Sweden. August 7–14, 1991; Amsterdam: Elsevier Science; 1994. p. 789-812.

Powell, JL. Estimation of semiparametric models. In: Engle, R.; McFadden, D., editors. Handbook of Econometrics. Vol. 4. Amsterdam: Elsevier; 1994. p. 2443-2521.

Reiersöl O. Confluence analysis by means of instrumental sets of variables. Arkiv för Matematik, Astronomi och Fysik. 1945; 32A(4):1–119.

Robins J. 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. 1986; 7(9–12):1393–1512.

Rosenbaum PR, Rubin DB. The central role of the propensity score in observational studies for causal effects. Biometrika. 1983 Apr; 70(1):41–55.

Rubin DB. Statistics and causal inference: Comment: Which ifs have causal answers. Journal of the American Statistical Association. 1986; 81(396):961–962.

Simon, HA. Causal ordering and identifiability. In: Hood, WC.; Koopmans, TC., editors. Studies in Econometric Method. Vol. Chapter 3. New York, NY: John Wiley & Sons, Inc.; 1953. p. 49-74.

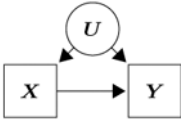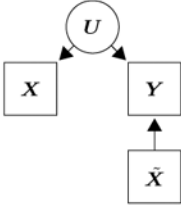Sobel ME. Discussion: 'the scientific model of causality'. Sociological Methodology. 2005; 35(1):99–133.

Spirtes, P. Proceedings of the Eleventh Conference Annual Conference on Uncertainty in Artificial Intelligence (UAI-95). San Francisco CA: Morgan Kaufmann; 1995. Directed cyclic graphical representations of feedback models; p. 491-498.

Spirtes, P.; Glymour, CN.; Scheines, R. Causation, Prediction and Search. 2. Cambridge, MA: MIT Press; 2000.

Tamer E. Incomplete simultaneous discrete response model with multiple equilibria. Review of Economic Studies. 2003 Jan; 70(1):147–165.

Tian, J.; Pearl, J. Proceedings of the Eighteenth National Conference on Artificial Intelligence. Cambridge, MA: AAAI Press; 2002. A general identification condition for causal effects; p. 567-573.

Tian, J.; Pearl, J. Technical report. Cognitive Systems Laboratory, University of California at Los Angeles; 2003. On the identification of causal effects.

Vytlacil EJ. Independence, monotonicity, and latent index models: An equivalence result. Econometrica. 2002 Jan; 70(1):331–341.

White H, Chalak K. Settable systems: An extension of Pearl's causal model with optimization, equilibrium, and learning. Journal of Machine Learning Research. 2009; 10:1759–1799.

Yule GU. On the correlation of total pauperism with proportion of out-relief. The Economic Journal. 1895; 5(20):603–611.

Step 1

Step 2

Step 3

$G$
Original
DAG

Possible
subgraphs

$G_{\overline{X}}$   $G_{\underline{X}}$

$G_?$

LMC

Conditional
Independence
Relationships

Step 6

Step 5

Step 4

Identification:
Find Empirical
Counterparts
(if any)

Implied
Causal
Parameters
(if any)

Apply
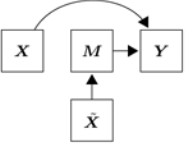Rule 1
Rule 2
Rule 3

**Figure 1. The Steps Required to Implement the Do-Calculus**

**Table 1**
**Haavelmo Empirical and Hypothetical Models**

| 1. Haavelmo Empirical Model | 2. Haavelmo Hypothetical Model |
|---|---|
| $\mathcal{T} = \{U, X, Y\}$<br>$\boldsymbol{\varepsilon} = \{\varepsilon_U, \varepsilon_X, \varepsilon_Y\}$<br>$Y = f_Y(X, U, \varepsilon_Y)$<br>$X = f_X(U, \varepsilon_X)$<br>$U = f_U(\varepsilon_U)$ | $\mathcal{T} = \{U, X, Y, \tilde{X}\}$<br>$\boldsymbol{\varepsilon} = \{\varepsilon_U, \varepsilon_X, \varepsilon_Y\}$<br>$Y = f_Y(X, \tilde{U}, \varepsilon_Y)$<br>$X = f_X(U, \varepsilon_X)$<br>$U = f_U(\varepsilon_U)$ |
|  |  |
| $Pa(U) = \emptyset,$<br>$Pa(X) = \{U\}$<br>$Pa(Y) = \{X, U\}$ | $Pa(U) = Pa(\tilde{X}) = \emptyset,$<br>$Pa(X) = \{U\}$<br>$Pa(Y) = \{X, \tilde{U}\}$ |
|  | $Y \perp\!\!\!\perp X \mid (X, \tilde{U})$<br>$X \perp\!\!\!\perp (\tilde{X}, Y) \mid U$<br>$X \perp\!\!\!\perp U$ |
| $Pr_E(Y, X, U) = Pr_E(Y\mid X, U)Pr_E(X\mid U)Pr_E(U)$ | $Pr_H(Y, X, U, \tilde{X}) = Pr_H(Y\mid X, \tilde{U})Pr_H(X\mid U)Pr_H(U)Pr_H(\tilde{X})$ |

This table has two columns and six panels separated by horizontal lines. Each column presents a causal model. The first panel names the models. The second panel presents the structural equations generating the models. In this row alone we make $\varepsilon$ explicit. In the other rows it is kept implicit to avoid clutter. Columns 1 and 2 are based on structural equations that have the same functional form, but have different inputs. The third panel represents the model as a DAG. Squares represent observed variables, circles represent unobserved variables. The fourth panel presents the parents in $\mathcal{T}$ of each variable. The fifth panel shows the conditional independence relationships generated by the application of the Local Markov Condition. The sixth panel presents the factorization of the joint distribution of variables in the Bayesian Network.

**Table 2**

**Models for Mediation Analysis**

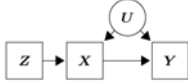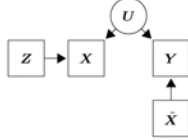| 1. Empirical Mediation Model | 2. Hypothetical Model for Indirect Effect of *X* on *Y* |
|---|---|
|  |  |

This table shows four models represented by DAGs. To simplify the displays we keep the unobservables in $\varepsilon$ implicit. Model 1 represents the empirical model for mediation analysis. The remaining three models are hypothetical models that target different causal effects of *X* on *Y*. Model 2 represents the hypothetical model for indirect effect *X* on *Y*.

**Table 3**

**Examples of the Do-calculus Notation**

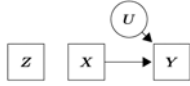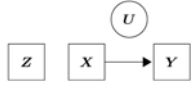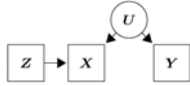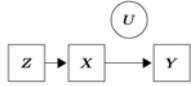| 1. Empirical Model *G* | 2. Modified Empirical Model $G_{X^-}$ |
|---|---|
|  |  |
| 3. Modified Empirical Model $G_{\overline{X},\underline{Z}}$ | 4. Modified Empirical Model $G_{\overline{X},\underline{Z(W)}}$ |
|  |  |

This table shows four models represented by DAGs that exemplify the notation of the do-calculus. Graph 1 is the original empirical model *G*. the remaining DAGs are generated by the deletion of arrows according to the do-calculus notation.

**Table 4**

**Instrumental Variable Empirical and Hypothetical Models**

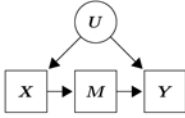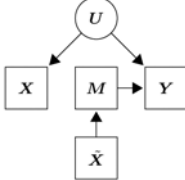| 1. Instrumental Variable Empirical Model | 2. Instrumental Variable Hypothetical Model |
|---|---|
| $\mathscr{T} = \{U, X, Z, Y\}$ <br> $\varepsilon = \{\varepsilon_U, \varepsilon_X, \varepsilon_Z, \varepsilon_Y\}$ <br> $Y = g_Y(X, U, \varepsilon_Y)$ <br> $X = g_X(U, Z, \varepsilon_X)$ <br> $Z = g_Z(\varepsilon_Z)$ <br> $U = g_U(\varepsilon_U)$ | $\tilde{\mathscr{T}} = \{U, X, Z, Y, \tilde{X}\}$ <br> $\varepsilon = \{\varepsilon_U, \varepsilon_X, \varepsilon_Z, \varepsilon_Y\}$ <br> $Y = g_Y(X, \tilde{U}, \varepsilon_Y)$ <br> $X = g_X(U, Z, \varepsilon_X)$ <br> $Z = g_Z(\varepsilon_Z)$ <br> $U = g_U(\varepsilon_U)$ |
|  |  |
| $Pa(U) = Pa(Z) = \varnothing,$ <br> $Pa(X) = \{U, Z\}$ <br> $Pa(Y) = \{U, X\}$ | $Pa(U) = Pa(U) = \varnothing,$ <br> $Pa(X) = \{U, Z\}$ <br> $Pa(Y) = \{U, \tilde{X}\}$ |
| $Z \perp\!\!\!\perp U$ <br> $Y \perp\!\!\!\perp Z\|(X, U)$ | $U \perp\!\!\!\perp (Z, \tilde{X})$ <br> $Z \perp\!\!\!\perp (U, X, \tilde{Y})$ <br> $X \perp\!\!\!\perp (U, Z, \tilde{X})$ <br> $X \perp\!\!\!\perp (X, \tilde{Y})\|(U, Z)$ <br> $Y \perp\!\!\!\perp (Z, X)\|(U, \tilde{X})$ |
| $Pr_E(Y, Z, X, U) = Pr_E(Y\|X, U)Pr_E(X\|U, Z)Pr_E(Z)Pr_E(U)$ | $Pr_H(Y, Z, X, U, \tilde{X}) = Pr_H(Y\|X, \tilde{U})Pr_H(X \| U, Z)Pr_H(Z)Pr_H(U)Pr_H(\tilde{X})$ |
| $Pr_E(Y, Z, U\|do(X) = x) = Pr_E(Y\|X = x, U)Pr_E(Z)Pr_E(U)$ | $Pr_H(Y, Z, X, U\|\tilde{X} = x) = Pr_H(Y\|\tilde{X} = x, U)Pr_H(X\|U, Z)Pr_H(Z)Pr_H(U)$ |

This table has two columns and seven panels separated by horizontal lines. Each column presents a causal model. The first panel names the model. The second panel presents the structural equations generating the model. In this row alone we make the $\varepsilon$ explicit. In the other rows it is kept implicit to avoid notational clutter. Columns 1 and 2 are based on structural equations that have the same functional form, but have different inputs. The third panel represents the model as a DAG. Squares represent observed variables and circles represent unobserved variables. The fourth panel presents the parents in $T$ of each variable. The fifth panel shows the conditional independence relationships generated by the application of the Local Markov Condition. The sixth panel presents the factorization of the joint distribution of variables in the Bayesian Network. The last panel of column 1 presents the joint distribution of variables when $X$ is fixed at $x$. ($do(X) = x$). The last panel of column 2 gives the joint distribution of variables generated by hypothetical models associated with empirical model 1 when $\tilde{X}$ is conditioned on $\tilde{X} = x$.

**Table 5**
**Do-calculus and the Instrumental Variable Model**

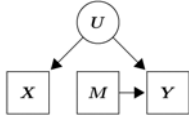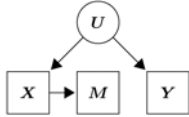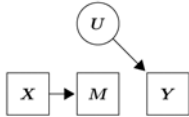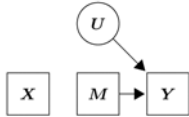| 1. Modified IV Model $G_{\underline{X}}$ | 2. Modified IV Model $G_{\overline{X},\underline{U}}$ |
|---|---|
|  |  |
| $U \perp\!\!\!\perp (Z, X)$ <br> $Z \perp\!\!\!\perp (U, X, Y)$ <br> $X \perp\!\!\!\perp (U, Z)$ <br> $Y \perp\!\!\!\perp Z\|(U, X)$ | $U \perp\!\!\!\perp (Z, X, Y)$ <br> $Z \perp\!\!\!\perp (U, X, Y)$ <br> $X \perp\!\!\!\perp (U, Z)$ <br> $Y \perp\!\!\!\perp (U, Z)\|X$ |
| 3. Modified IV Model $G_{\overline{U},\underline{X}}$ | 4. Modified IV Model $G_{\underline{U}}$ |
|  |  |
| $U \perp\!\!\!\perp Z$ <br> $Z \perp\!\!\!\perp (U, Y)$ <br> $X \perp\!\!\!\perp Y\|(U, Z)$ <br> $Y \perp\!\!\!\perp (Z, X)\|U$ | $U \perp\!\!\!\perp (Z, X, Y)$ <br> $Z \perp\!\!\!\perp Z$ <br> $X \perp\!\!\!\perp U\|Z$ <br> $Y \perp\!\!\!\perp (U, Z)\|X$ |

This table shows four models represented by DAGs ($\varepsilon$ are kept implicit to avoid notational clutter). Squares represent observed variables, circles represent unobserved variables. Each DAG is generated by the deletion of arrows of the original Instrumental Variable model (first column of Table 4) according to the rules of the do-calculus. Below each model, we show conditional independent relations generated by the application of the Local Markov Condition (5) to variables of the models.

**Table 6**
**"Front-Door" Empirical and Hypothetical Models**

| 1. Pearl's "Front-Door" Empirical Model | 2. Our Version of the "Front-Door" Hypothetical Model |
|---|---|
| $\mathcal{T} = \{U, X, M, Y\}$ <br> $\varepsilon = \{\varepsilon_U, \varepsilon_X, \varepsilon_M, \varepsilon_y\}$ <br> $Y = f_Y(M, U, \varepsilon_Y)$ <br> $X = f_X(U, \varepsilon_X)$ <br> $M = f_M(X, \varepsilon_M)$ <br> $U = f_U(\varepsilon_U)$ | $\mathcal{T} = \{U, X, M, Y, \tilde{X}\}$ <br> $\varepsilon = \{\varepsilon_U, \varepsilon_X, \varepsilon_M, \varepsilon_Y\}$ <br> $Y = f_Y(M, U, \varepsilon_Y)$ <br> $X = f_X(U, \varepsilon_X)$ <br> $M = f_M(\tilde{X}, \varepsilon_M)$ <br> $U = f_U(\varepsilon_U)$ |
|  |  |
| $Pa(U) = \emptyset$ <br> $Pa(X) = \{U\}$ <br> $Pa(M) = \{X\}$ <br> $Pa(Y) = \{M, U\}$ | $Pa(U) = Pa(\tilde{X}) = \emptyset$ <br> $Pa(X) = \{U\}$ <br> $Pa(M) = \{\tilde{X}\}$ <br> $Pa(Y) = \{M, U\}$ |
| $Y \perp\!\!\!\perp X \vert (M, U)$ <br> $M \perp\!\!\!\perp U \vert X$ | $Y \perp\!\!\!\perp (X, \tilde{X}) \vert (M, U)$ <br> $M \perp\!\!\!\perp (U, X) \vert \tilde{X}$ <br> $X \perp\!\!\!\perp (M, \tilde{X}, Y) \vert U$ <br> $U \perp\!\!\!\perp (M, \tilde{X})$ <br> $\tilde{X} \perp\!\!\!\perp (X, U)$ |
| $Pr_E(Y, M, X, U) = Pr_E(Y\vert M, U)Pr_E(X\vert U)Pr_E(M\vert X)Pr_E(U)$ | $Pr_H(Y, M, X, U, \tilde{X}) = Pr_H(Y\vert M, U) \, Pr_E(X\vert U)Pr_H(M\vert\tilde{X})Pr_H(U)Pr_H(\tilde{X})$ |
| $Pr_E(Y, M, U\vert do(X) = x) = Pr_E(Y\vert M, U)Pr_E(M\vert X = x)Pr_E(U)$ | $Pr_H(Y, M, U, X\vert\tilde{X} = x) = Pr_H(Y\vert M, U) \, Pr_E(X\vert U) \, Pr_H(M\vert\tilde{X} = x)Pr_H(U)$ |

This table has two columns and seven panels separated by horizontal lines. Each column presents a causal model. The first panel names the models. The second panel presents the structural equations generating the model. In this row alone we make $\varepsilon$ explicit. In the other it is kept implicit to avoid clutter. Columns 1 and 2 are based on structural equations that have the same functional form, but have different inputs. The third panel represents the model as a DAG. Squares represent observed variables, circles represent unobserved variables. The fourth panel presents the parents in $\mathcal{T}$ of each variable. The fifth panel shows the conditional independence relationships generated by the application of the Local Markov Condition. The sixth panel presents the factorization of the joint distribution of variables in the Bayesian Network. The last panel of column 1 presents the joint distribution of variables when $X$ is fixed at $x$ using the "do operator." The last panel of column 2 gives the joint distribution of variables generated by the hypothetical models associated with empirical model 1 when $\tilde{X}$ is conditioned at $\tilde{X} = x$.

**Table 7**

**Do-calculus and the Front-Door Model**

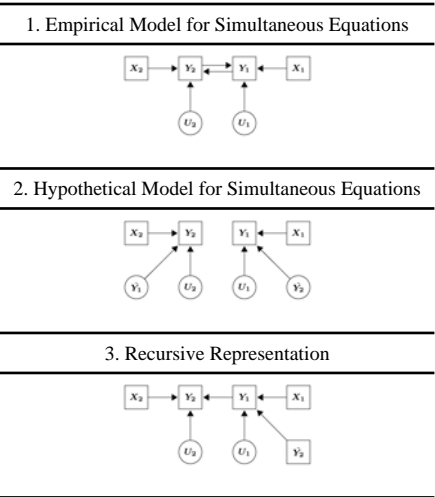| 1. Modified Front-Door Model $G_{\underline{X}} = G_{\overline{\underline{M}}}$ | 2. Modified Front-Door Model $G_{\overline{M}}$ |
|---|---|
|  |  |
| $(Y, M) \perp\!\!\!\perp X \mid U$<br>$(X, U) \perp\!\!\!\perp M$ | $(X, M) \perp\!\!\!\perp Y \mid U$<br>$(Y, U) \perp\!\!\!\perp M \mid X$ |
| 3. Modified Front-Door Model $G_{\overline{X}, \underline{M}}$ | 4. Modified Front-Door Model $G_{X, \tilde{M}}$ |
|  |  |
| $(X, M) \perp\!\!\!\perp (Y, U)$ | $(Y, M, U) \perp\!\!\!\perp X$<br>$U \perp\!\!\!\perp M$ |

This table shows four models represented by DAGs ($\varepsilon$ are kept implicit to avoid notational clutter). Squares represent observed variables, circles represent unobserved variables. Each DAG is generated by the deletion of arrows of the original Front-Door model (first column of Table 6) according to the rules of the do-calculus. Below each model, we show conditional independent relations generated by the application of the Local Markov Condition (5) to variables of the models.

**Table 8**

**Summarizing the Do-calculus of Pearl (2009) and the Haavelmo Approach**

| **Common Features** of Haavelmo and Do-Calculus: | | |
|---|---|---|
| **Autonomy** (Frisch, 1938) | | |
| **Errors Terms:** $\varepsilon$ mutually independent | | |
| **Statistical Tools:** LMC and Graphoid Axioms apply | | |
| **Counterfactuals:** Fixing or Do-operator is a causal, not statistical, operation | | |
| **Distinctive Features** of Haavelmo and Do-Calculus: | | |
| | **Haavelmo** | **Do-calculus** |
| **Approach:** | Thinks outside the box of the empirical model by constructing a new hypothetical model motivated by, but distinct from, the empirical model where fixing can be analyzed using standard tools of probability | Thinks inside the box of the empirical model; Creates complex nonstandard rules to introduce fixing into a probabilistic framework |
| **Introduces:** | Constructs a hypothetical model | Graphical rules |
| **Identification:** | Connects $Pr_H$ and $Pr_E$ | Iteration of do-calculus rules |
| **Versatility:** | Basic statistical principles apply | Creates new rules of statistics |

**Table 9**

**Models for Simultaneous Equations**

---

1. Empirical Model for Simultaneous Equations



---

2. Hypothetical Model for Simultaneous Equations



---

3. Recursive Representation



---

This table shows two models. (The variables in $\varepsilon$ are kept implicit.) Model 1 represents the empirical model for Simultaneous Equations where $Y_1$ and $Y_2$ cause each other. Model 1 is cyclic, and hence it is not a DAG. Model 2 represents one possible hypothetical model associated with the empirical model for Simultaneous Equations. In Model 2, the hypothetical variable $\tilde{Y}_2$ is associated with the causal link of $Y_2$ on $Y_1$ of Model 1 and the hypothetical variable $\tilde{Y}_1$ is associated with the causal link of $Y_1$ on $Y_2$ of Model 1. Model 3 presents a recursive ("causal chain") representation of a hypothetical model in which simultaneity is broken in the original model and the resulting $Y_1$ is set outside of the model through hypothetical variation. This thought experiment is one way to produce the model on the left hand side of panel 2.