

Reflections on the 2021 Nobel Memorial Prize Awarded to David Card, Joshua Angrist, and Guido Imbens

LENNART B. ACKERMANS
Erasmus University Rotterdam

The 2021 *Sveriges Riksbank Prize in Economic Sciences in Memory of Alfred Nobel* was awarded in two halves. One half was awarded to David Card “for his empirical contributions to labour economics” (The Royal Swedish Academy of Sciences 2021, 1). The other half was awarded jointly to Joshua D. Angrist and Guido W. Imbens “for their methodological contributions to the analysis of causal relationships” (1). In this article, I (a philosopher of science interested in causal inference in economics) reflect on the second half of the 2021 Nobel Prize, awarded to Angrist and Imbens.

Two beautiful examples of causal inference in economics are Angrist (1990) and Angrist and Krueger (1991), published shortly after Joshua Angrist obtained his PhD in 1989. (His co-laureates David Card and Guido Imbens are his contemporaries, obtaining their PhDs in 1983 and 1991, respectively.) The 1990 study estimates the causal effect of veteran status on earnings 30 years later. It finds that white U.S. veterans from the Vietnam War have approximately 15% lower earnings as a result of military service. The 1991 study finds that, in the U.S., having an additional year of compulsory schooling has a large effect on earnings later in life (Angrist and Krueger 1991).

At the time, other studies into both of these subjects struggled to find causal effects, as opposed to mere correlations. For example, when one finds a negative correlation between veteran status and earnings, it is unclear whether this is because veterans had a lower earning potential prior to being enlisted, or because lower earnings result from serving in the military. Angrist (1990) solved this problem using an *instrumental variable approach*. During the Vietnam War, men were drafted based on a lottery that assigned numbers between 1 and 365 based on birth date. Only men below a certain lottery number were drafted. Since the draft lottery

number is randomly assigned, the causal effect of the lottery number on earnings can be identified from the observed data. Since an observational situation like this is similar to an experiment, such as a randomised controlled trial (RCT), it is called a *natural experiment*.

However, we are ultimately not interested in the effect of lottery number on earnings, but in the effect of military service. Lottery numbers are not a perfect substitute for enlistment because some people with low lottery numbers did not enlist and some people with high numbers volunteered. To estimate the effect of military service itself, Angrist used the lottery number as an *instrumental variable*. (An instrumental variable has some effect on the explanatory variable of interest. It is used for its better inference properties—particularly *unconfoundedness*, discussed in section 2 below.) An instrumental variable approach first estimates the effect of lottery number on earnings and the effect of lottery number on military enlistment. With some additional assumptions that might be controversial, the two results can be combined to derive the effect of military service on earnings.

The draft lottery study has become a classic example of successful causal inference, but in 1990 Angrist was not yet convinced of its persuasiveness. As he recalls in his Prize Lecture: “Guido and I soon began asking each other: What really do we learn from the *draft eligibility* and *quarter of birth* natural experiments?” (Angrist 2021). It was only in the years that followed that Angrist and Imbens were able to answer this question with their ground-breaking methodological work, which includes classic papers such as “Identification and Estimation of Local Average Treatment Effects” (Imbens and Angrist 1994) and “Identification of Causal Effects Using Instrumental Variables” (Angrist, Imbens, and Rubin 1996). Later work added additional causal methods such as *regression discontinuity* approaches (Imbens and Lemieux 2008). It is in methodological studies such as these where the most important contributions of Angrist and Imbens lie.

An uncontroversial way to characterise these contributions is to say that Angrist and Imbens developed clever methods of causal inference—such as instrumental variable approaches and regression discontinuity—which have allowed economists to produce more successful causal studies, such as the draft lottery and compulsory schooling studies. However, I shall argue that their contribution is more important. In the 1980s, econometrics was in a state of crisis. Edward Leamer put it well in his 1983

article titled, “Let’s Take the Con Out of Econometrics”, in which he pronounced that econometric practice of the time was “decidedly unscientific” (37). Even worse, everyone knew it: “Hardly anyone takes data analyses seriously. Or perhaps more accurately, hardly anyone takes anyone else’s data analyses seriously” (37). (See also LaLonde 1986 and discussions of the credibility crisis by Angrist and Pischke 2010; Imbens 2022.) Angrist and Imbens led the way out of this crisis by shifting the field’s attention toward causal research design. Three decades later, Joshua Angrist and Jörn-Steffen Pischke (2010) declared that econometrics had made significant progress since Leamer’s critique.

I agree. Advances in causal methodology made by Angrist, Imbens, and others have been especially important because they mark the first step out of the dark ages of econometrics. However, this process is still incomplete, and the legacy of Angrist and Imbens will grow larger still—or at least, it *should* grow larger. This brings me to my second thesis. At present, the econometric methodology of the crisis era is still prevalent. Moreover, economists are divided on how to proceed, with heated debates over the causal framework to adopt. The framework championed by Angrist and Imbens, the *Rubin Causal Model*, has limitations that I believe might hamper a more widespread adoption in the field. I will argue that the profession needs to prioritise the resolution of these problems so that it can put causal inference at the forefront of economic inquiry.

I. THE IDENTIFICATION PROBLEM

While it may not have been clear in 1991, the draft lottery and compulsory schooling studies were important achievements that marked the way out of the credibility crisis. But why was there a crisis in the first place, why were these instrumental variable studies so successful, and how did they contribute to resolving the crisis? In order to answer these questions, let me introduce the *identification problem*.

In simple terms, the identification problem is that statistics is not causal inference. More precisely, exercises in *pure statistics* are never enough to discover a causal relation. (There are a variety of other definitions of the identification problem in economics. Sometimes it refers to the more general problem that a parameter in a model cannot be estimated from observations.)

By *pure statistics* I mean mathematically describing data in the following way. It is assumed that the data in one’s sample is drawn from a larger set of data called the population—which can be hypothetical and infinitely

large—either randomly or using some known or unknown procedure. The pure statistician uses the sample data to make inferences about the population data. Typically, one is interested in a number of mathematical parameters describing the population data, such as means, conditional means, and correlation coefficients. Often, the statistical parameters of interest figure as Greek letters in regression equations describing the population data, such as:

$$Y_i = \alpha A_i + \beta B_i + \epsilon_i. \quad (1)$$

Here (Y_i, A_i, B_i) is the i 'th data point and α and β are population parameters that the statistician aims to estimate from the sample data, typically using a regression method such as ordinary least squares (OLS).

The error term ϵ_i is the deviation of Y_i from its expected value given (A_i, B_i) , which is estimated by the regression residual. Equations such as (1) can be great for describing data statistically. In that case, they may also be valuable for predicting the values of new samples drawn from the same population.

An entirely different interpretation of equation (1) is that it is a model that describes an underlying causal structure, in which case it is called a *structural equation*, *structural model*, or *causal model*. If (1) is a causal model, α is interpreted as a causal effect of A on Y , β as a causal effect of B on Y , and ϵ_i as a combination of unobserved variables having a separate causal effect on Y . These are now *causal parameters* instead of statistical parameters. On this interpretation, the equation describes how the i 'th individual's outcome Y_i was causally determined based on the values of A_i , B_i , and ϵ_i . It implies that individual i 's value of Y_i would have been as described by (1) if i 's values of A_i , B_i , and ϵ_i had been given any different value. Such counterfactual knowledge importantly allows the researcher to predict what will happen if the studied population changes, for example, due to policy interventions.

If (1) is a causal model, the researcher can still try to estimate the (causal) parameters with the methods of pure statistics, such as OLS. However, such an endeavour is usually unwarranted. A major problem is that the same data can always be generated by multiple different mechanisms. Even if the data is well described by (1), there are many different models that could have generated the exact same data. If the pure statistician would have estimated the parameters of another model, entirely

different causal parameters would have been ‘discovered’. Hence, statistical estimates of α and β can be interpreted as causal effects *only if* we have good reasons to assume that the underlying causal structure is given by (1).

This suffices to show that causal inference requires more than pure statistics. That said, it should be mentioned as an aside that the use of advanced statistical methods is extremely important for causal inference. Some important contributions that Angrist and Imbens have made to causal methodology are best classified as pure statistics. For example, Angrist, Imbens, and Krueger (1999) offer a solution to a problem of bias that occurs in two-stage-least-squares (2SLS) estimation, a statistical method essential for causal methods using instrumental variables. This solution, called *Jackknife Instrumental Variables Estimation* (JIVE), is then used to re-analyse the 1991 Angrist and Krueger study.

There have been economists throughout the 20th century who understood the identification problem and the difficulties with causal inference (e.g., Haavelmo 1943). Nevertheless, econometrics textbooks to this day—while great at teaching pure statistics—are creating more confusion than clarity when it comes to causal inference, as several authors in the field now recognise (Heckman and Pinto 2022b; Angrist, Imbens, and Rubin 1996).¹

When a typical textbook in econometrics introduces the OLS regression model, it informs the reader of a crucial assumption (e.g., Wooldridge 2010; Greene 2018). The error term ϵ_i , as it appears for example in equation (1), must be uncorrelated with the other variables appearing on the right-hand side. This assumption is known as *econometric exogeneity* (or just *exogeneity*). If the error term does correlate with a regressor, called *econometric endogeneity*, then the regression estimates of the parameters α and β are said to be biased or non-causal. Once exogeneity is assumed, econometricians are able to apply powerful tools from pure statistics. However, exogeneity is an assumption about structure (causal or otherwise) and needs a defence that goes beyond pure statistics.

This is where the textbooks fall short. The exogeneity assumption as it typically appears in works of econometrics is “either meaningless or false”, as Pratt and Schlaifer (1984, 11) summarise it. Econometricians tend to define error terms as the combined effect of unobserved variables. However, the concept of ‘omitted variables’—without giving it a more precise definition—is so vague that the exogeneity assumption has no real

¹ See footnote 4 for the relevant quotes from these authors.

content. In any given regression an infinite number of variables have been omitted. Without specifying *which* unobserved variables are meant, the correlation between ϵ_i and the regressors is not defined. Depending on how ϵ_i is interpreted, the correlation with other variables can have any value.

To see why this is so, consider the model $Y = \alpha X + \epsilon_1$, where X is exogenous, i.e., $\text{Cov}(X, \epsilon_1) = 0$ and X and ϵ_1 have mean 0. Now consider another variable, $\epsilon_2 = \beta X + \gamma \epsilon_1$. Then the model $Y = \alpha X + \epsilon_2$ could describe the data just as well. However, in this model, X correlates with ϵ_2 . Without specifying what a variable ϵ is, it could be ϵ_1 , ϵ_2 , or many other things. Which variable one should choose depends on the causal effect one wants to measure. For example, if the causal effect of X on Y is α , then the error term is ϵ_2 , and X is endogenous—so the causal effect cannot be identified with OLS. If the causal effect is β , then X is exogenous and the causal effect identifiable.

Hence, for the exogeneity assumption as it is typically invoked in econometric studies, it is impossible to check whether it is true or false. That said, if one assumes a particular causal structure, it is possible to give the error term a definition for which its correlation with other variables is defined (Pearl 2009, chapter 5). However, textbook econometrics is devoid of such causal assumptions. Under these circumstances, exogeneity is not a meaningful assumption. (These problems have long been understood. See the classic papers by Haavelmo 1943 and Pratt and Schlaifer 1984.)

From my understanding, what was wrong with econometrics as practised in the 1980s was that researchers did not have a clear understanding of the above issues (this diagnosis is similar to Pearl 2009, chapter 5; Imbens 2022). In particular, the difference between pure statistics and causal inference was often obscured, with regression equations like (1) not having a clear interpretation as either describing data or causal structure. As Imbens (2022) observes, the term ‘causality’ was rarely used in econometrics between the 1960s and 1980s, until it was revived in the 1990s—despite the fact that econometricians were often concerned with clearly causal questions.

This lack of causal terminology can even be found in the very articles that identified the credibility crisis in the 1980s. Leamer (1983), while seemingly concerned that conventional regression estimates in econometrics do not match causal parameters, does not mention causality in the

paper. Another interesting paper is LaLonde (1986), which compares experimental and non-experimental methods using the same data. The data comes from an RCT designed to estimate the effect on trainee earnings of an employment program. Putting aside the control group, LaLonde applied state-of-the-art econometric techniques for use with observational data—and he was unable to replicate the results from the RCT. Like Leamer, LaLonde did not use words like ‘causality’. However, LaLonde did identify ‘model misspecification’ in observational methods as a problem. If a ‘model’ is a causal model, this was going in the right direction. However, to properly analyse and resolve the problems that econometricians became aware of in the 1980s, a more principled understanding of causation and methods of causal inference was needed.

Given the clarity of causal reasoning found in the draft lottery and compulsory schooling papers (Angrist 1990; Angrist and Krueger 1991), I imagine that the authors had a clear understanding of the difficulties associated with causal inference. The semi-experimental methods they used were uncommon in economics at the time and would later lead the way out of the credibility crisis. However, the causal reasoning in these papers is not principled in the sense of being based on well-studied formal principles of causal inference. Around 1990, causal reasoning in economics relied on intuitions rather than theory and was thus more of an art than a science. But this was about to change.

II. THE RUBIN CAUSAL MODEL

Causal inference, like statistics, must be done with the help of formal frameworks that assist the scientist in reasoning correctly and precisely. The problem in the 1980s was that economists had mastered well-developed and sophisticated tools of pure statistics, while their tools for causal inference were lagging behind. Fortunately, statisticians had already developed a framework for causal inference, known as the Rubin Causal Model (Rubin 1974) named after Donald Rubin by Holland (1986), but going back to Neyman ([1923] 1990) and Cox (1958). This section introduces the Rubin Causal Model (RCM) and illustrates how it improved econometricians’ understanding of causal methods, using the example of instrumental variables.

The strategy of RCM is to use the RCT as a foundation on which to build a framework which extends well beyond RCTs. As in an experiment, we imagine that each individual can be given the treatment ($T = 1$) or no treatment ($T = 0$). An individual’s outcome if treated is denoted $Y_i(1)$, and

an individual's outcome if not treated is denoted $Y_i(0)$. These are called *potential outcomes*, of which at least one is counterfactual. The *individual treatment effect* for i is given by $Y_i(1) - Y_i(0)$.

It is a virtue of the RCM that it relates a causal effect so clearly to a counterfactual: the effect of i 's treatment is the difference between i 's outcome if i were treated and if i were not treated. Unfortunately, only one of these outcomes can be observed. Hence, we need clever strategies in order to learn something about causal effects without ever being able to observe them directly.

As it turns out, various types of *average* treatment effects (ATEs) can sometimes be derived from statistical data. This is the case, for example, for an RCT with perfect compliance. (Perfect compliance means that all participants get the treatment if and only if they are assigned the treatment). In the perfect RCT, due to the random assignment of treatment, the average observed difference between the treatment and control group is an estimate of the average counterfactual difference for all individuals. More precisely, one can show that

$$E[Y_i(1) - Y_i(0)] = E[Y_i | T = 1] - E[Y_i | T = 0],$$

where Y_i is i 's observed outcome. The expectation on the left is called the *average treatment effect*—which is an average of causal effects that is not directly observable. The expression on the right, on the other hand, can be estimated from the observed data with the techniques of pure statistics.

Unfortunately, observational data is never like an RCT with perfect compliance. But fortunately, there is now a large literature with methods to identify treatment effects with weaker assumptions, including assumptions that are sometimes satisfied in observational data. The important contributions from Angrist and Imbens lie mostly in this area.

Their most influential achievement is perhaps a method for identifying the *local average treatment effect* or *LATE* (Imbens and Angrist 1994; Angrist, Imbens, and Rubin 1996). The LATE is an average treatment effect for a subpopulation of the data: namely, those individuals whose treatment status always matches their assignment, called *compliers*. The LATE can be estimated in RCTs with imperfect compliance, but its greatest success stems from the fact that it can sometimes be estimated from

purely observational data. This is the case when the data contains an instrumental variable—call it Z_i —which has the properties of treatment assignment in an imperfect RCT.

To illustrate the kind of assumptions required for causal inference within RCM, I will give a somewhat technical discussion of the LATE using RCM terminology. This will pay off in the next section, which compares the merits and problems of RCM with other frameworks.

An instrumental variable Z_i is an observable variable that has some causal influence on individuals' treatment $T_i(z)$ with $z \in \{0,1\}$. Here $T_i(z)$ is the treatment that i would have if it were the case that $Z_i = z$. The LATE is defined as

$$E[Y_i(1) - Y_i(0) \mid T_i(1) = 1, T_i(0) = 0].$$

Angrist and Imbens showed that the LATE can be identified if three important assumptions are satisfied (as well as some others). First, the potential outcomes $Y_i(t)$ are unaffected by Z_i . More precisely, if $Y_i(z, t)$ is i 's potential outcome given $Z_i = z$ and $T_i = t$, then we have $Y_i(z, t) = Y_i(t)$, for all $z, t \in \{0,1\}$. (A more intuitive formulation of this assumption may be that Y_i is unaffected by Z_i if the treatment T_i is held fixed). This is called the *exclusion restriction*. Second, Z_i must have the properties of random assignment. In RCM terminology, this assumption states that Z_i is probabilistically independent of the potential outcomes $(Y_i(0), Y_i(1), T_i(0), T_i(1))$. (That is, it is jointly independent of these four variables. I will explore this assumption in greater detail in the next section.) This assumption is usually called *unconfoundedness*. Third, assignment to the treatment must make treatment more likely for each individual. More precisely, there should be no *defiers*, individuals who do the opposite of their treatment assignments. Defiers are individuals such that $T_i(1) = 0$ and $T_i(0) = 1$.

The LATE method showcases how RCM can be used to prove mathematically that a causal effect can be identified from the data given these assumptions. This subsequently makes it possible for applied researchers to increase the credibility of their studies, provided that they can make it plausible that these assumptions are indeed satisfied. The assumptions contained in LATE and other RCM-based methods are certainly easier to defend than econometric exogeneity, by virtue of their rigorous explication. However, they are still not *quite* easy to defend—which brings me to one of RCM's foremost shortcomings (see also Pearl 2009, 98–102).

Let us look at what these assumptions mean in the draft lottery study. In this study, the instrumental variable Z_i is $Z_i = 1$ if the individual has a low lottery number, such that he is eligible for the draft, $Z_i = 0$ otherwise. For the treatment we have $T_i(Z_i) = 1$ if the individual is enlisted, $T_i(Z_i) = 0$ if not. Finally, $Y_i = Y_i(T_i)$ is i 's observed income 30 years after the draft.

The exclusion and no-defiers assumptions are relatively straightforward to defend. The no-defiers assumption says that there are no individuals that would have volunteered for military service with a high lottery number but would not enlist with a low lottery number. It seems safe to assume that such individuals are rare enough to ignore.

On the other hand, assessing unconfoundedness is a mental nightmare. Unconfoundedness says that lottery number Z_i is jointly independent of all the potential outcomes $(Y_i(0), Y_i(1), T_i(0), T_i(1))$. Assessing this assumption requires one to imagine population data which contains not only individuals' actual treatment and outcome values but also the treatment and outcomes that they would have under counterfactual conditions. Without additional guidance, this assumption is very hard to assess. Unfortunately, the causal framework RCM itself does not provide much help in assessing whether unconfoundedness is satisfied.

Methods to test indirectly whether unconfoundedness is satisfied, based on RCM, do exist (see e.g., Imbens and Rubin 2015, chapter 21). The problem is that assessing unconfoundedness requires much more than some mathematical methods which a researcher can simply 'run'. More importantly, it requires an informal understanding of the underlying causal structure and a way to *translate* this understanding into formal assumptions of probabilistic independencies. (Even the tests in Imbens and Rubin 2015 require informal input based on the researcher's intuitions and theoretical knowledge.)

It is my position that a causal framework is supposed to help the researcher with this translation step from structural causal knowledge to methodological assumptions. RCM, however, is unsuited for this task by construction. Causal connections are not expressed in RCM, which instead focuses on independencies in imaginary population data that includes potential outcomes. The result is that all assumptions in RCM are expressed in terms of imaginary data, without using any causal terms. To assess the assumptions, however, one needs to consult one's causal knowledge. For example, knowledge of whether Z_i has common causes with T_i or Y_i should be used to assess unconfoundedness. Such knowledge comes in terms of causal connections, not imaginary population data.

However, proponents of RCM insist that they find these assumptions quite intuitive. For example, in response to a similar concern voiced by Pearl, Imbens (2020, 1164) replies: “I think that statement [from Pearl] misses the point. This setting, where the critical assumption is ignorability or unconfoundedness, is so common and well studied that merely referring to its label is probably sufficient for researchers to understand what is being assumed”.

Irrespective of its potential problems, RCM has been extremely important for the development of causal methods such as the LATE. Both Angrist and Imbens mention this importance in their prize lectures (Angrist 2021; Imbens 2021). Nevertheless, RCM has not managed to replace the textbook approach to econometrics in most econometric research. Part of the problem is that there are several contenders aiming to replace textbook econometrics as a framework for causal reasoning.

III. CONTENDING CAUSAL FRAMEWORKS

Separately from the RCM developed by statisticians, computer scientists and philosophers developed another causal framework, which I will call the *Pearl Causal Model* (PCM) after its primary author Judea Pearl (Pearl and Verma 1991; Spirtes, Glymour, and Scheines 1993; Pearl 1995, 2009). The PCM makes extensive use of directed acyclic graphs (DAGs) to formulate assumptions about causal structure. James Heckman has defended another causal framework which he claims is closer to the traditional econometric framework (Heckman 2000, 2005; Heckman and Pinto 2015). Let us call this the *Heckman Causal Model* (HCM). In this section, I summarise these frameworks and show how they can be used to shed light on the assumption of unconfoundedness.

Both Heckman and Pearl are influenced by earlier economists’ work on structural equation modelling such as Frisch ([1938] 1995) and Haavelmo (1943, 1944). Moreover, in its most recent explication, HCM makes heavy use of DAGs to express structural causal assumptions graphically, as well as other tools from the PCM literature (Heckman and Pinto 2015, 2022b). Hence, the two approaches are spiritually and practically similar. RCM, on the other hand, eschews the use of structural equations.

In both approaches, the foundation of causal inference is causal modelling. Before one can reliably estimate causal effect sizes, one typically needs to have knowledge about causal structure—that is, knowledge

about how the variables in a system are causally connected. Causal models in HCM and PCM summarise such information using equations and graphs. For example, figure 1 represents the research design of an instrumental variable setup like the draft lottery study. The nodes in this graph represent causal variables, and the arcs represent causal connections. For example, $Z_i \rightarrow T_i$ means that Z_i is a cause of T_i . The unobserved variable U_i is responsible for individual differences in their response to the treatment assignment, as represented by the arc $U_i \rightarrow Y_i$.

HCM and PCM can express unconfoundedness both graphically and probabilistically. In what follows I illustrate how these frameworks give the researchers additional tools in understanding and assessing unconfoundedness.

Consider again the model of an instrumental variable setup in figure 1. The graphical unconfoundedness assumption states that Z_i is causally connected with Y_i only by causing it (via T_i , as in figure 1). In other words, there must not be a common cause C_i of Z_i and Y_i , that is, a path

$$Z_i \leftarrow C_i \rightarrow Y_i.$$

Graphical unconfoundedness implies probabilistic unconfoundedness under an assumption called the *Causal Markov Condition*. The Causal Markov Condition states that a causal variable is independent of its non-descendants conditional on its parents, supposing that the DAG is a sufficiently accurate representation of reality. For instance, in figure 1, the Causal Markov Condition implies that Y_i is independent of Z_i given (T_i, U_i) . The Causal Markov Condition is a well-studied principle that is plausible in most circumstances, although objections exist (e.g., Cartwright 1999).

Based on the Causal Markov Condition and a variety of rules for manipulating conditional independence relations (from Dawid 1979), a researcher can quickly derive all conditional independencies implied by a DAG. To make a connection with the probabilistic unconfoundedness assumption in the previous section, one can create ‘hypothetical versions’ of a graph in which treatment variables are replaced by counterfactual variables. In HCM, one creates a hypothetical model given counterfactual assignments of Z_i as follows. First, add a counterfactual treatment assignment variable \hat{Z}_i to the graph. Then remove all outgoing arrows from Z_i and instead assign them as outgoing arrows from \hat{Z}_i . The resulting graph, depicted in figure 2, represents the causal model given counterfactual assignments \hat{Z}_i . (PCM uses a slightly different procedure to create counterfactual models.) By reading off independencies from the hypothetical

graph, the researcher can quickly observe that Z_i is independent of T_i given counterfactual assignments \hat{Z}_i . Similarly, figure 3 gives the hypothetical graph given counterfactual assignments of T_i . From this graph, the researcher can observe that Z_i is independent of Y_i given counterfactual assignments of \hat{T}_i . These results in turn can be shown to imply the probabilistic unconfoundedness assumption from the previous section.²

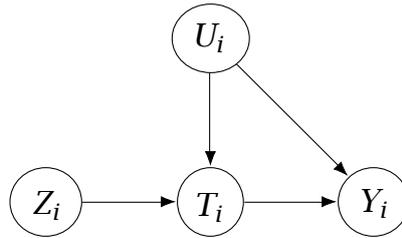


Figure 1: Causal graph of an instrumental variable setup.

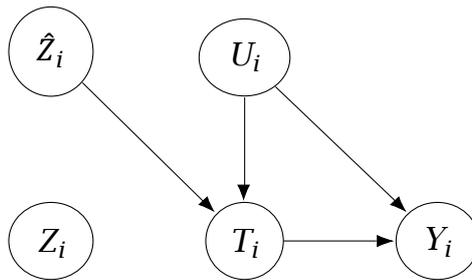


Figure 2: Causal graph for counterfactual assignments of Z_i .

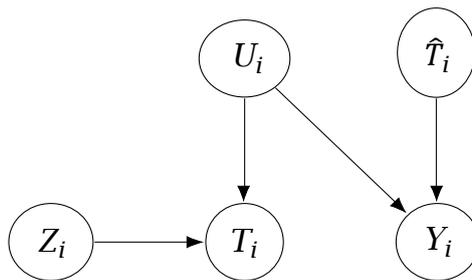


Figure 3: Causal graph for counterfactual assignments of T_i .

The above illustrates how causal graphs can be used by researchers to use their theoretical knowledge of causal structure, as expressed in a DAG, to assess whether assumptions for causal methods are satisfied.

² See Heckman and Pinto (2022b) for a more detailed analysis.

While above I went through the reasoning from structure to independence assumptions explicitly, researchers do not typically need to do so themselves. Pearl's book, and Heckman and Pinto's recent articles, describe many causal methods based on assumptions that are expressed in graphical terms, such as Pearl's *back-door criterion* and *front-door criterion* (Pearl 2009). This allows researchers to immediately apply these methods once they have identified an accurate causal structure.

The graphical approach of PCM and HCM shifts the researcher's attention to an important precondition of causal inference: the identification of causal structure. As illustrated above, one needs detailed knowledge of causal structure before the assumptions required for causal methods can be verified. Graphical causal frameworks not only make it easier to express structural causal knowledge but also come with a rich literature that helps researchers to discover causal structure from data, including algorithms that search for causal relations in data (Spirtes, Glymour, and Scheines 2001).

Hence, the graphical frameworks PCM and HCM supply the researchers with a more complete set of tools, including tools for estimating causal effect sizes, verifying structural assumptions, and discovering causal structure. All tools are part of the same graphical framework, allowing scientists to combine them easily.

IV. THE RECENT DEBATE: HECKMAN AND IMBENS

The previous section showcases some of the benefits of adopting a graphical approach to causal inference. Given these benefits, there is a good case to be made that economists should adopt HCM or PCM instead of RCM. However, not everyone agrees. Both Angrist and Imbens are vocal proponents of RCM and have written textbooks that exclusively rely on RCM (Angrist and Pischke 2009; Imbens and Rubin 2015). In an article published a year before he won the Nobel Prize, Guido Imbens criticises PCM (and indirectly, HCM), claiming that RCM is better suitable for empirical practice in economics—while acknowledging that the graphical approach “has not had as much impact in economics as it should have” (2020, 1130). Given the similarities between HCM and PCM, many of Imbens' criticisms of PCM apply to HCM as well. On the other hand, Heckman and Pinto (2022b) argue that HCM is a more suitable framework for economists than both RCM and PCM. This section summarises the debate and argues that the latest advances in Heckman's camp give the graphical approach an edge over the others.

The 2020 article by Guido Imbens is a great overview of the arguments in support of RCM. First, contrary to what I have argued above, Imbens claims that the formulation of key assumptions is, in fact, *more* intuitive in RCM than in graphical frameworks. According to Imbens, the RCM formulations “capture the way researchers think of causal relationships” (2020, 1130). Second, RCM is claimed to connect more easily to traditional economic models such as the supply and demand model. Interestingly, Heckman and Pinto make the exact opposite claim, arguing that RCM as well as PCM “have significant limitations when applied to the wide variety of problems that economists face” (2022b, 894). Third, while Imbens acknowledges that PCM is advantageous for complex models with many variables, he claims that such models “are not particularly popular in empirical economics” (2020, 1155). Fourth, RCM is useful for dealing with the problem of treatment effect heterogeneity. Fifth, RCM is claimed to connect better with many practical questions of causal study design and the inference of causal effects.

The sixth and most forceful reason for preferring RCM (in my opinion) is that it is better capable of capturing the assumptions required for some causal methods. By reasoning about probabilistic independencies directly—bypassing considerations of structure—RCM has undoubtedly allowed methodologists to discover methods that would otherwise be overlooked because they seem improbable if you have a graphical perspective. Instrumental variable methods—the LATE in particular—are an example of this. With causal structures as in figure 1, the effect of T_i on Y_i can only be identified given additional non-graphical assumptions such as the no-defiers assumption. This is recognised by the others in the debate as well (Pearl 2009, 90; Heckman and Pinto 2022b, 913).

However, James Heckman and Rodrigo Pinto’s (2022b) recent work demonstrates that HCM is in fact extremely versatile. It is capable of formulating the assumptions needed for instrumental variable methods such as LATE, as well as those needed for methods from PCM, such as front-door and back-door adjustment. Hence, the most apparent advantage of RCM, that it has a natural way of explicating assumptions needed for instrumental variable methods, may no longer be a relative advantage compared to HCM. At the same time, HCM has all the advantages of PCM by virtue of incorporating graphical models, as I illustrated in the previous section.

V. CONCLUSION: ALL ECONOMETRICIANS SHOULD ADOPT CAUSAL FRAMEWORKS

Hence, based on the most recent developments, it seems to me that HCM has an edge over the other frameworks. It is versatile, suitable for many empirical methods in economics, and deeply rooted in economic tradition. However, Heckman and Pinto may go a bit too far when they say that the use of RCM and PCM by economists has been detrimental:

Many econometricians and applied economists now emulate what they read in statistics or computer science journals. They have forgotten or never learned their own field's foundational work to the detriment of rigorous causal policy analysis.³ (Heckman and Pinto 2022a)

The above claim is somewhat misleading, given the serious problems with the econometric approach as taught in textbooks and still practised today. This tradition is *responsible* for the problems in econometrics that became apparent in the 1980s. Causal frameworks such as RCM, on the other hand, have greatly contributed to the development of sound causal methods in econometrics. Heckman and Pinto may mean that the field's founders from which they draw inspiration, such as Haavelmo and Frisch, had a better (and causal) understanding of structural equation models than what is found in textbooks. They are right about that, but this older tradition was forgotten or corrupted in the later 20th century (see Pearl 2009, section 5.1.2). Moreover, authors within PCM and RCM *also* claim to be inspired by Haavelmo's work. It may be more accurate to say that all present-day causal frameworks draw on early 20th-century work, while none of the current causal frameworks can claim to stand in a continuous tradition from then until the present.

Both the RCM and HCM sides of the debate now seem to agree that the textbook definition of econometric exogeneity is inadequate, preferring alternative concepts from the newer causal frameworks.⁴ While Heckman's earlier causal framework still relies on econometric exogeneity (Heckman 2005), Heckman and Pinto's recent version no longer makes

³ The published version Heckman and Pinto (2022b) makes a similar point.

⁴ Angrist, Imbens, and Rubin say about exogeneity: "Typically the researcher does not have a firm idea what these disturbances [error terms] really represent, and therefore it is difficult to draw realistic conclusions or communicate results based on their properties" (1996, 446). Imbens calls econometric exogeneity "inadequate" (1997, 93). Heckman and Pinto say: "Econometrics textbooks often discuss causality as a property of an estimator, usually ordinary least squares (OLS). This approach reverses the logic of causality. It also generates confusion, since the OLS model is described by statistical assumptions that are void of causality" (2022b, 896).

any references to econometric exogeneity (Heckman and Pinto 2015, 2022b). This is a clear way in which all frameworks depart from the econometric tradition. It is also a good thing. Although econometric exogeneity when defined precisely and in structural terms can be a helpful concept—as argued by Pearl (2009, 169–170)—economists on both sides are abandoning the ambiguous textbook definition of econometric exogeneity and replacing it with clearly defined causal assumptions.

All of the causal frameworks on offer are a significant improvement to the field of econometrics. What is troubling, however, is that the textbook approach to econometrics is largely unchanged. The typical econometrics textbook has as its foundation the OLS regression model and the econometric exogeneity assumption, while RCM might be discussed much later as an afterthought (see e.g., Wooldridge 2010; Greene 2018). These textbooks have one important improvement compared to earlier days: they recognise that causal identification is the fundamental problem that economists are concerned with. For example, the first sentence of the introduction in Wooldridge reads: “The goal of most empirical studies in economics and other social sciences is to determine whether a change in one variable, say w , causes a change in another variable, say y ” (2010, 3). Hence, it is surprising that these textbooks take an approach that the leading experts on causal inference in the field—including Heckman, Imbens, and Angrist—recognise as inadequate.

The textbook approach has consequences for econometric practice. For the world’s star economists, causal frameworks might not be absolutely essential. After all, Angrist was able to produce interesting and credible causal studies in the early 1990s without relying on RCM. However, he was doing so at a time in which econometric research was widely believed to be incredible by its own practitioners. Causal frameworks are essential for the standardisation of credible causal methods and for bringing these methods to a larger group of researchers.

Moreover, it can be shown that practising economists make mistakes as a direct result of the confusion created by the concept of econometric exogeneity, as I do in Ackermans (2022, appendix A). There I discuss a complicated type of sensitivity analysis invented by economists to estimate the size of causal bias. However, the method is incapable of improving the estimate of causal bias already assumed by the researcher’s choice of parameters. These kinds of useless mathematical exercises can be avoided if modern causal frameworks are used as the foundation of training and practice in econometrics.

Why is progress in econometric education so slow? Perhaps the field simply needs more time. But one factor must be that economists cannot agree on which causal framework should be adopted. Without a consensus on this matter, textbook authors have little incentive to overthrow the approach they have taken for decades and which is currently used more widely than any of the modern causal frameworks.

Like many in the debate, I have strong views on the respective merits of the different frameworks. However, what is more important than *which* causal framework to adopt is that *a* causal framework is adopted—since PCM, HCM, and RCM are all big improvements over textbook econometrics. The profession should resolve the dispute about causal frameworks and update its graduate teaching. That is the only way to solidify the advances in causal methodology made by Angrist, Imbens, and others and assist future generations of economists in further advancing their work.

REFERENCES

- Ackermans, Lennart B. 2022. "Causal bias in measures of inequality of opportunity." *Synthese* 200 (6): 429.
- Angrist, J. D., G. W. Imbens, and A. B. Krueger. 1999. "Jackknife instrumental variables estimation." *Journal of Applied Econometrics* 14 (1): 57-67.
- Angrist, Joshua D. 1990. "Lifetime Earnings and the Vietnam Era Draft Lottery: Evidence from Social Security Administrative Records." *The American Economic Review* 80 (3): 313-336.
- Angrist, Joshua D. 2021. "Prize Lecture." NobelPrize.org.
- Angrist, Joshua D., Guido W. Imbens, and Donald B. Rubin. 1996. "Identification of Causal Effects Using Instrumental Variables." *Journal of the American Statistical Association* 91 (434): 444-455.
- Angrist, Joshua D., and Alan B. Krueger. 1991. "Does Compulsory School Attendance Affect Schooling and Earnings?" *The Quarterly Journal of Economics* 106, no. 4 (November): 979-1014.
- Angrist, Joshua D., and Jörn-Steffen Pischke. 2009. *Mostly Harmless Econometrics: An Empiricist's Companion*. Princeton, NJ: Princeton University Press.
- Angrist, Joshua D., and Jörn-Steffen Pischke. 2010. "The Credibility Revolution in Empirical Economics: How Better Research Design Is Taking the Con out of Econometrics." *Journal of Economic Perspectives* 24, no. 2 (June): 3-30.
- Cartwright, Nancy. 1999. "Causal Diversity and the Markov Condition." *Synthese* 121 (1): 3-27.
- Cox, D.R. 1958. *Planning of Experiments*. New York, NY: Wiley.
- Dawid, A. P. 1979. "Conditional Independence in Statistical Theory." *Journal of the Royal Statistical Society: Series B* 41 (1): 1-15.
- Frisch, Ragnar. (1938) 1995. "Autonomy of economic relations: Statistical versus theoretical relations in economic macrodynamics." In *The Foundations of Econometric Analysis*, edited by David F. Hendry and Mary S. Morgan, 407-425. Cambridge: Cambridge University Press.

- Greene, William H. 2018. *Econometric Analysis*. 8th ed. New York, NY: Pearson.
- Haavelmo, Trygve. 1943. "The Statistical Implications of a System of Simultaneous Equations." *Econometrica* 11 (1): 1-12.
- Haavelmo, Trygve. 1944. "The Probability Approach in Econometrics." *Econometrica* 12 (Supplement): iii-115.
- Heckman, James J. 2000. "Causal Parameters and Policy Analysis in Economics: A Twentieth Century Retrospective." *The Quarterly Journal of Economics* 115, no. 1 (February): 45-97.
- Heckman, James J. 2005. "The Scientific Model of Causality." *Sociological Methodology* 35 (1): 1-97
- Heckman, James J., and Rodrigo Pinto. 2015. "Causal Analysis After Haavelmo." *Econometric Theory* 31 (1): 115-151.
- Heckman, James J., and Rodrigo Pinto. 2022a. *Causality and Econometrics*. Working Paper 29787. National Bureau of Economic Research, February.
- Heckman, James J., and Rodrigo Pinto. 2022b. "The Econometric Model for Causal Policy Analysis." *Annual Review of Economics* 14 (1): 893-923.
- Holland, Paul W. 1986. "Statistics and Causal Inference." *Journal of the American Statistical Association* 81 (396): 945-960.
- Imbens, Guido W. 1997. Review of *The Foundations of Econometric Analysis* by David F. Hendry and Mary S. Morgan. *Journal of Applied Econometrics* 12 (1): 91-94.
- Imbens, Guido W. 2020. "Potential Outcome and Directed Acyclic Graph Approaches to Causality: Relevance for Empirical Practice in Economics." *Journal of Economic Literature* 58, no. 4 (December): 1129-1179.
- Imbens, Guido W. 2021. "Prize Lecture." NobelPrize.org.
- Imbens, Guido W. 2022. "Causality in Econometrics: Choice vs Chance." *Econometrica* 90 (6): 2541-2566.
- Imbens, Guido W., and Joshua D. Angrist. 1994. "Identification and Estimation of Local Average Treatment Effects." *Econometrica* 62 (2): 467-475.
- Imbens, Guido W., and Thomas Lemieux. 2008. "Regression discontinuity designs: A guide to practice." The regression discontinuity design: Theory and applications, *Journal of Econometrics* 142 (2): 615-635.
- Imbens, Guido W., and Donald B. Rubin. 2015. *Causal Inference for Statistics, Social, and Biomedical Sciences*. New York, NY: Cambridge University Press.
- LaLonde, Robert J. 1986. "Evaluating the Econometric Evaluations of Training Programs with Experimental Data." *The American Economic Review* 76 (4): 604-620.
- Leamer, Edward E. 1983. "Let's Take the Con Out of Econometrics." *The American Economic Review* 73 (1): 31-43.
- Neyman, Jerzy. (1923) 1990. "On the Application of Probability Theory to Agricultural Experiments. Essay on Principles. Section 9." Translated by D. M. Dabrowska and T. P. Speed. *Statistical Science* 5 (4): 465-472.
- Pearl, Judea. 1995. "Causal diagrams for empirical research." *Biometrika* 82, no. 4 (December): 669-688.
- Pearl, Judea. 2009. *Causality*. 2nd ed. Cambridge: Cambridge University Press.
- Pearl, Judea, and T.S. Verma. 1991. "A Theory of Inferred Causation." In *Proceedings of the Second International Conference on Principles of Knowledge Representation and Reasoning*, 441-452. Cambridge, MA, USA: Morgan Kaufmann Publishers Inc.

- Pratt, J. W., and Robert Schlaifer. 1984. "On the Nature and Discovery of Structure." *Journal of the American Statistical Association* 79 (385): 9-21.
- The Royal Swedish Academy of Sciences. 2021. "The Prize in Economic Sciences 2021". Press release, October 11, 2021.
- Rubin, Donald B. 1974. "Estimating causal effects of treatments in randomized and non-randomized studies." *Journal of educational Psychology* 66 (5): 688-701.
- Spirtes, Peter, Clark Glymour, and Richard Scheines. 1993. *Causation, Prediction, and Search*. New York, NY: Springer-Verlag.
- Spirtes, Peter, Clark Glymour, and Richard Scheines. 2001. *Causation, Prediction, and Search*. 2nd ed. Cambridge, MA: MIT Press.
- Wooldridge, Jeffrey M. 2010. *Econometric Analysis of Cross Section and Panel Data*. 2nd ed. Cambridge, MA: MIT Press.

Lennart B. Ackermans is a PhD candidate at Erasmus University Rotterdam. He works on philosophy of science, epistemology, and ethics.
Contact e-mail: <ackermans@esphil.eur.nl>