

THE ROLE OF HYPOTHESIS TESTING  
IN THE  
MOLDING OF ECONOMETRIC MODELS

BY KEVIN D. HOOVER

CHOPE Working Paper No. 2012-03

January 2012



CENTER FOR THE  
HISTORY OF POLITICAL ECONOMY  
AT DUKE UNIVERSITY

**The Role of Hypothesis Testing  
in the  
Molding of Econometric Models**

Kevin D. Hoover  
Department of Economics and Department of Philosophy  
Duke University  
Box 90097  
Durham, North Carolina 27708-0097

Tel. (919) 660-1876  
E-mail [kd.hoover@duke.edu](mailto:kd.hoover@duke.edu)

Keynote lecture from the Tilburg-Madrid Conference on Hypothesis Tests: Foundations and Applications, Universidad Nacional de Educación a Distancia (UNED) Madrid, Spain, 15-16 December 2011. The support of the U.S. National Science Foundation (grant no. NSF SES-1026983) is gratefully acknowledged.

Revised, 18 January 2012

Abstract

of

**The Role of Hypothesis Testing in the Molding of Econometric Models**

The paper is a keynote lecture from the Tilburg-Madrid Conference on Hypothesis Tests: Foundations and Applications at the Universidad Nacional de Educación a Distancia (UNED) Madrid, Spain, 15-16 December 2011. It addresses the role of tests of statistical hypotheses (specification tests) in selection of a statistically admissible model in which to evaluate economic hypotheses. The issue is formulated in the context of recent philosophical accounts on the nature of models and related to some results in the literature on specification search.

**Keywords:** statistical testing, hypothesis tests, models, general-to-specific specification search, optional stopping, severe tests, costs of search, costs of inference, extreme-bounds analysis, LSE econometric methodology

**JEL Codes:** B41, C18, C12, C50

## I. Introduction

When David Teira invited me to give this talk, he first suggested that I might represent a frequentist view of statistics, while other people would represent a Bayesian view. He rightly saw that my prejudices ran in the frequentist direction. But as I told him immediately, even though I have long been interested in the philosophical aspects of econometrics, the aspects that have attracted most of my interest are not in the foundations of probability. In particular, I have avoided saying much about Bayesianism, feeling that I did not actually know enough to have a well-grounded opinion: “better to keep silent and be thought a fool than to open your mouth and remove all doubt.”<sup>1</sup> I have learned a lot more about Bayesianism since that first demurrer, and I am sure that I cannot avoid at least implicitly taking a stand on that debate; still, I am not qualified to serve as a principal in such a controversy.

Later David suggested to me that I might reflect on the controversy over significance testing in economics in which my colleague Mark Siegler and I were engaged with Deirdre McCloskey and Stephen Ziliak (Hoover and Siegler 2008a, b; McCloskey and Ziliak 1996; Ziliak and McCloskey 2004). In a series of papers, McCloskey and Ziliak presented the use of significance tests as a major failure of economics and economists as dunderheaded incompetents when it comes to quantitative methods. You will get the spirit of the controversy from McCloskey’s (2002) hyperbolic rhetoric:

The progress of science has been seriously damaged. You can’t believe anything that comes out of [tests of statistical significance]. Not a word. It is all nonsense, which future generations of economists are going to have to do all over again. Most of what appears in the best journals of economics is unscientific rubbish. [p. 55]

...

---

<sup>1</sup> Often attributed in one version or another to Mark Twain.

Until economics stops believing . . . that an intellectual free lunch is to be gotten from . . . statistical significance . . . our understanding of the economic world will continue to be crippled by the spreading, ramifying, hideous sins. [p. 57]

The message of our side of the debate was: “*Take comfort, things are not so dark as all that* (Hoover and Siegler 2008a, p. 31). McCloskey and Ziliak had grabbed onto a bone – that statistical significance and economic (or substantive) significance were not the same thing – and like pit bulls they were ready to maul all in their path to hold on to it. Our principal objections were to their slipshod scholarship, including their unwarranted calumny against econometrics and the economics profession, and to their incoherent methodological position. Of course, we did not dispute – and indeed we know of no economist (and they gave not a single example of one) who in principle does dispute – the validity of that distinction. But what then?

The point of distinguishing substantive and statistical significance is developed with vastly more insight in the work of Deborah Mayo and Aris Spanos on “severe testing” (Mayo 1996; Mayo and Spanos 2006). The idea of severe testing may be unfamiliar, so let me give a précis. Merely accepting or rejecting a null hypothesis is a weak test on Mayo and Spanos’s view. Consider a test of a null hypothesis. One accepts the null if the test statistic is less than the critical value for a designated size and rejects it if is greater. The  $p$ -value gives the *attained size* – that is, that test size that the actual data would have supported. But is such a test severe? That depends on the alternative hypothesis. We must choose an alternative that is just big enough to matter. A test is severe if, under that alternative, the probability of the opposite conclusion reached in our test of the null is highly probable. Thus, a test is severe if we give it every chance to fail and yet it still succeeds. Severity is clearly related to the power of a test in Neyman and Pearson’s sense; but, by analogy with the attained size measured by the  $p$ -value, it is measured by *attained power*, measured not at a prespecified critical value but at the actual value of the test statistic derived from the observed data. A reasonable conclusion would be that

economists and other empiricists should pay attention to the severity of their tests, but not ban them.

Rather than revisit the controversy with McCloskey and Ziliak, which reached methodologically diminishing returns long ago, I would prefer to say something about an issue to which I have contributed positively and in which the idea of severe testing is at home.

Economics is a modeling science. The Nobel laureate James Heckman (2000, p. 46) has said that, just as the Jews are the “people of the book,” the economists are “the people of the model.” Of course, economists are not alone in this. In the period since the mid-20<sup>th</sup> century, the model has become the dominant epistemic tool in a wide variety of sciences. The philosophy of science used to pay a great deal of attention to issues such as the axiomatic structure of formal scientific theories and to demarcation criteria between science and nonscience. Increasingly, it has focused on how models work in science.<sup>2</sup> The change is part of the “naturalistic turn” in the philosophy of science – the laudable notion that, if we want to know how science works, we ought to try to understand the practices of scientists.

My own work as a methodologist derives from my work as a monetary and macroeconomist. I was interested in the role of monetary and monetary policy in controlling inflation and real output, which raises questions about the causal structure of the economy. In that context, I developed a kind of interventionist or “natural-experiments” approach to causal inference (Hoover 2001). In implementing the approach in real-world cases, I was forced to characterize the data statistically and adopted the model-selection strategies of the LSE (London School of Economics) approach of David Hendry and his colleagues and coworkers (Mizon 1984, 1995; Hendry 1987, 2000). In the event, it was this approach, which relies heavily on statistical testing, and my own contribution to causal inference that raised questions with

---

<sup>2</sup> See Morrison and Morgan (1999) and Morgan (2012).

referees. So, I came to the problems of statistical testing through the backdoor. Even now, I prefer to keep my reflections grounded in the specific problems encountered in my own practices.

Issues related to statistical testing can, I think, be subsumed to more general issues related to modeling. Frequently, inferential problems assume that the form of a probability distribution is known and the test relates to some parameters of that distribution. McCloskey and Ziliak provide a neat example of what worries me: “the accuracy of [the] estimated mean [of a regression coefficient] depends on the properties of the error term, the specification of the model, and so forth. *But to fix ideas suppose that all the usual econometric problems have been solved*” (McCloskey and Ziliak 1996, p. 98, emphasis added). They, like many others, ignore the larger problem: how would we justify such the supposition that all the usual econometric problems have been solved? All my own work on causality in macroeconomics was about choosing the form of the relationships that McCloskey and Ziliak and most econometric textbooks simply take as given. I want to consider the role of statistical tests in addressing the problem of selecting – or better, *shaping* or *molding* – economic models.<sup>3</sup>

## II. Models

Let us begin with models, without supposing that they are necessarily stochastic or invoke probability. The concept of causation, on my preferred account, is one of mechanism or structure (Hoover 2001, chs. 1-4). The object of an empirical analysis of causation is to construct a model that recapitulates the salient features of the mechanism and displays its causal architecture perspicaciously. In most cases, the role of a model is to make hidden causal relationships visible. Economic data do not wear their causal relationships on their faces. But

---

<sup>3</sup> I am echoing here Boumans (2005, chs. 1 and 3 ) notion of the “mathematical moulding of economic theory.”

that is a matter of degree. Some modeling exercises recapitulate relationships that are, as it were, visible to the naked eye. For example, children and aeronautical engineers make models of airplanes in which the mapping from the real airplane to the model is not much of a mystery. (I do not wish to underplay the complexity of the relationship of the modeled to the model even in this case; see Sterrett 2005)

Models are instruments for relating truths about the world. Although models are sometimes "approximations" in an exact sense of that word, I prefer to think of them, up to some explicit or implicit level of precision, as telling the plain truth about limited aspects of the world or from particular perspectives on the world (Hoover, forthcoming). Models may have varying levels of precision and cast the world from various points of view, but their premier virtue is accuracy (i.e., in being used to claim what in fact happens in the world).

Models are governed by the constitutive properties of the model. Some of these properties are specific to the model and irrelevant to the world. A wind-tunnel model, for example, need not have an internal structure that mimics an actual airplane, so long as the mass and exterior shapes are appropriate. Models are also governed by internal rules of operation. They may be closed systems in which deductive results are available or their operation may be only analogical with results available through simulation. In either case, the world of the model is not automatically informative about the real world. It will be informative only if there is a good mapping between model and world on relevant dimensions, which adds an interpretive relationship between model as object and its implications for the real world.

Perhaps *the* principal function of models is as engines for counterfactual analysis. We validate the mapping between real-world and the model using observations of the real world as our guide, but the utility of the model is that manipulation reveals facts about the world that we



have not yet or, perhaps, cannot ever observe directly. This is the source of the utility of a model for prediction or control.<sup>4</sup>

The general characteristics of models are evident in such transparent cases as the model airplane. In economics, however, as in many disciplines, we value models precisely because they clarify the actions of hidden mechanisms. Consider Project Ultra in which the British successfully read German military codes in World War II. They constructed a working model of the German’s Enigma code machine. In part, they benefitted from stealing versions of the machine. Nevertheless, a substantial part of their success arose from figuring out how the machine must work in order to generate the intercepted coded signals. Their model did not need to be an exact copy; it did need to be an appropriate analogue. And it served as a tool of counterfactual analysis: given that the model provides a mechanism that accounts for some observed code with a particular initial setting, the machine allowed the code-breakers to determine what any particular piece of plain text would look like with some other initial setting.

The process of modeling the Enigma machine was not a process of conjecture and refutation or of hypothesis testing of the form, “propose a hypothesis and then ask, ‘accept or reject?’” Rather it was a process of molding the model mechanism to constraints – some directly from data, some from other considerations. And it is a process very unlike the philosophers’ accounts of inductive logic. Typically, induction is presented as a problem of moving from instances to a generalization: Raven<sub>1</sub> is black, raven<sub>2</sub> is black, raven<sub>3</sub> is black, . . . raven<sub>n</sub> is black; therefore, all ravens are black (or very probably all ravens are black). The form of inference is not true to typical scientific and practical reasoning, which is much closer to the reasoning of the code-breakers. First, it is too simple. It may be a good strategy for finding the proportion of white beans in an urn, but it fails to come to grips with the wide range of inferential

---

<sup>4</sup> See Hoover (2011) for the role of economic models in counterfactual analysis.

patterns found in science and everyday life. Second, it does not deal with the role of creativity in learning. We really must engage in a good deal of guessing of answers on the basis of preexisting beliefs. This process is not, however, unfettered. It is a process in which our beliefs are mutually constraining, even when those beliefs are not held with complete conviction. We gain conviction from their mutual reinforcement.

Here is a mundane illustration. In many cases when we have solved a complex crossword puzzle, our conviction that our solution is correct is nearly absolute. It is not that there could not be a world in which an entirely different set of answers fit the physical constraints of the puzzle grid and satisfied reasonable interpretations of the clues. We cannot rule that out *a priori*, but neither need we let it have great force in our thinking when the fact is that our solution all fits together nicely and that it is extremely difficult to get a solution to fit together at all. In solving the puzzle, we have passed a severe test.

Creative imagination is essential to forward progress, but the limits of imagination also constrain the alternative choices that we might consider. Frequently, the imaginations of different investigators point to different solutions, which must be checked against the commonly accepted constraints or tested by generating new constraints that may not satisfy one or other alternative. We can, for example, see the Ptolemaic and Copernican models of the solar system as different imaginary solutions to the observed motions of the planets and stars. Our preference for Copernicus over Ptolemy is that ultimately, though this was not immediately obvious, it better fit the constraints. Of course, the original Copernican system is not entirely satisfactory, and our modern model has been molded to adapt to the additional constraints of later observations and our belief in Newton's laws, among other things.

### III. Economic Models

The problem of empirical economics is largely one of inferring the nature and properties of the hidden mechanisms of the economy. We do that in the manner of the code-breakers: we construct analogue models of some features of the economy. Economic theory can be regarded as a set of model templates for such mechanisms, and the problem of the applied economist is to find a good template and to mold it to various constraints imposed by observed data and preexisting beliefs.

Let me give a hackneyed example. Suppose that we want to know how the price of electricity affects the demand for electricity. We might appeal to a supply-and-demand model:

(1)	$Q_E^D = a + bP_E + cT$	Demand
(2)	$Q_E^S = d + eP_E + fP_C$	Supply
(3)	$Q = Q_E^S = Q_E^D$	Equilibrium

where  $Q$  = quantity of electricity;  $Q_E^D$  = demand for electricity;  $Q_E^S$  = supply of electricity;  $P_E$  = price of electricity;  $P_C$  = price of coal; and  $T$  = temperature. Figure 1 shows the model in a graphical form. Here a problem is evident: if we know only the data ( $Q$ ,  $P_E$ ,  $P_C$ , and  $T$ ), we account for a single observation where the supply and demand curves cross and we cannot learn what we want to learn, namely how price affects the demand for electricity. If it happened that  $T$  were constant and  $P_C$  were variable and, in addition, some other assumptions held, then shifts in the supply curve (shown in the figure as grey lines) would trace out the demand curve and we would be able to *identify* the values for the coefficients  $a$  and  $b$ . If both  $T$  and  $P_C$  varied, then we would be able to identify all of the coefficients.

But what about those other assumptions? They include that variations in  $T$  and  $P_C$  are independent of each other, that the underlying relationships are well modeled as linear, that  $T$

does not appear in the supply equation nor  $P_C$  in the demand equation, that there are no other shifters of the equations, and so forth. That knowledge is not in the observable data. How do we know it? The standard answer to this question among economists, going back at least to Haavelmo's seminal "Probability Approach in Econometrics" (1944), is that it is *a priori* knowledge based in economic theory. But how did we come to have such knowledge? Indeed, this question is hardly ever addressed. The concept of *a priori* knowledge, which is relied upon to do a vast amount of work, has never to my knowledge been examined by econometricians or economic methodologists. And the professed faith in economic theory as the source of such knowledge amounts to whistling in the dark.<sup>5</sup> Economic theory in its pure form generates very weak conclusions: for example, we can reasonably hold it to suggest that demand curves slope down ( $b < 0$ ), but it certainly does not tell us that demand depends on temperature ( $T$ ) and not on the price of coal ( $P_C$ ) or any other factor. Sometimes we are told that it is not theory, but subject-matter knowledge (expert knowledge) that supplies the ground for our *a priori* knowledge. This is, perhaps, closer to the truth, but equally unanalyzed by econometricians, methodologists, and philosophers alike. To answer a question about the nature of demand, we need to have a model with known properties that maps well onto properties of the world. Is there a systematic method for obtaining such knowledge? Would the statistical methods used in econometrics help? The answer must be, no, if econometrics, as it is presented in many (perhaps most) textbooks, is limited to the problem of statistical estimation of the parameters of structures assumed to be known in advance.

---

<sup>5</sup> Skepticism about identification does have a Liu (1960) and Sims (1980); cf. Hoover (2006).

#### **IV. Econometric Models**

The problem of *a priori* knowledge and of identification are typically thought of as econometric or statistical problems. The supply-and-demand model shows, however, that the problem arises in deterministic systems. It is a problem of modeling and not a problem of probability or statistics *per se*. The problem is to find sufficient constraints to allow us convincingly to mold our model into one that is strongly analogous to the hidden mechanisms of the economy. We cannot do that by armchair speculation or appeals to weak economic theory. The only hope is for the data to provide some of the key constraints in the same manner as they do in solving a crossword puzzle or breaking a code. If the world is indeterministic, either ontologically (reality is deeply stochastic) or epistemically (we are so ignorant of the full spectrum of causes that from our limited point of view reality acts as if it were deeply stochastic), we will need to account for its indeterminism in our models. We may do this by developing probabilistic models. (There may, of course, be other modeling tools applicable to indeterministic models. We are too apt to privilege our analytical creations. There is no more reason to assume that well-known treatments of probability provide the only possible resource for confronting indeterminism than there is for thinking that balsa wood is the only suitable material for model airplanes.)

Probability is not a property of the world; it is a property of our model. And probability models grab on to the world in just the same way as other models do through analogy in specific respects useful for particular purposes of particular agents (cf. Giere 2006, p. 60). Giere treats models as predicates, such as "is red." For example, a classical particle system is a model of behavior that obeys the Newton's laws and the law of gravity for interacting point masses. To say that our solar system is a classical particle system is to make a claim that this model provides accurate analogies for the motions of the planets around the sun (Giere 1979, ch. 5; also Giere

1999, pp. 98-100, 122 and 2006, p. 65; cf. Hausman 1992, p. 74). Kolmogorov's (or other) axiomatizations of probability provide just such a model of probability and can be regarded as a predicate in the same manner. The cases that most interest me are cases where the laws of probability can be accurately predicated of processes in the economy or the physical world. A model can be predicated wherever it effectively captures analogous features; so I by no means rule out predicating the probability model of beliefs either descriptively or normatively.

Statistical tests comes into modeling on this view as measures of the aptness of the predication. A cooked example, originally due to Johansen 2006, pp. 293-295) will help to make my point (also see Hoover, Johansen, and Juselius 2008, pp. 252-253). Johansen starts with the unobservable data-generating process:

$$(4) \quad \begin{aligned} x_t &= 0.9x_{t-1} + 1.0 + \varepsilon_t; \\ t &= 1, 2, \dots, 100; \quad x_0 = 10, \end{aligned}$$

where the  $\varepsilon_t$  are identically independently distributed (i.i.d)  $N(0,1)$ . Note that  $E(x_t) = 1/(1 - 0.9) = 10$  and  $\text{var}(x_t) = 1/(1 - 0.9^2)$ . Consider an economic theory that predicts that the mean value of  $x$  is  $\mu = 10$ . (Here, the theory happens to be *exactly* true, but it need not always be so.) To test the theory we need to provide a model of the probability process. One model is:

### Model 1

$$(5) \quad x_t = \mu + v_t,$$

where the  $v_t$  are i.i.d.  $N(0, \sigma^2)$ . For one simulation of equation (4), the maximum-likelihood estimate Model 1 yields  $\hat{\sigma}^2 = 5.256$  and an estimate of a 95 percent asymptotic confidence interval for  $\hat{\mu}$ :  $\hat{\mu} \pm 1.96\hat{\sigma}/\sqrt{T} = 9.138 \pm 0.449$ . Since 10 does not lie within the confidence interval, it might appear, then, that we have good grounds to reject the hypothesis that  $\mu = 10$ .

But is this model accurately predicated of the data-generating process? The error terms in Model 1 are i.i.d normal. Given the data-generating process (1), a simple statistical test would most likely show that the residuals do not conform to that assumption, but are serially correlated. We can conclude, then, that Model 1 cannot be accurately predicated of the data-generating process and that our estimate  $\mu$  is unlikely to be properly analogous to  $E(x_t)$ , which is its target and, consequently, our theory has not been tested adequately tested.

An alternative statistical model is

### Model 2

$$(6) \quad x_t = \rho x_{t-1} + \mu(1 - \rho) + v_t,$$

where again the  $v_t$  are i.i.d.  $N(0, \sigma^2)$  and  $E(x_t) = \mu$ , if  $|\rho| < 1$ . Model 1 is nested in Model 2. Again omitting details, the maximum likelihood estimate of Model 2 for the same simulated data yields an estimates of  $\hat{\rho} = 0.923$  and  $\hat{\sigma}^2 = 0.744$ , which translates into the 95 percent asymptotic confidence interval of  $\hat{\mu} \pm 1.96\hat{\sigma} / [(1 - \rho)\sqrt{T}] = 9.123 \pm 2.247$ . On the basis of this confidence interval, we cannot reject  $\mu = 10$ .

Statistical tests play two different roles in Johansen's cooked illustration. First, they translate the data into constraints on the form of the model in the same way that the puzzle grid and reasonable interpretations of the clues impose constraints on the solution to the crossword puzzle. Model 1 does not display serially correlation ( $\rho = 0$ ). It is highly unlikely that a model of that form could generate the pattern of the observed data, so we conclude that it would be inaccurate to predicate Model 1 of the data-generating process. Model 2 allows us to compare the estimated  $\hat{\rho}$  to a null of  $\hat{\rho} = 0$ . The test rejects the null, and relative to an alternative such as  $\rho = 0.9$ , the test is severe in the sense of Mayo and Spanos. The way in which Model 1 fails actually suggests a

property that any more accurate model will have – i.e., it must be able to generate serially correlated realizations.

The second role of statistical tests in the cooked illustration is the more familiar one: they are used to evaluate hypotheses conditional on the form of the model. If Model 2 is an acceptable model, then  $\mu$  is not very precisely estimated, but it is consistent with the hypothesis that  $\mu = 10$ . This is the basis on which hypothesis testing is usually conducted. The model is given, and we are concerned entirely with the precision of the estimates.

To interpret an estimate of a parameter, we must have a model in which the parameter is meaningful. Econometricians are wont to say that economic theory provides that model. While economic theory may impose some constraints on acceptable models, it is a vanishingly small class of cases in which it provides a single, estimable model. The first use of statistical models is to draw on the resources of the data itself to cover the weakness of economic theory in this regard. Seen this way, the first use of statistical tests in molding the model shows that Model 1 is not an acceptable starting place for the second use of statistical tests. The precision of the estimate of  $\mu$  is spurious, because that estimate takes its meaning from a model that does not accurately analogize to a salient feature of the world.

Econometrics as it is taught in textbooks – and even as it is sometimes practiced – focuses on the second use of statistical tests as if we had *a priori* knowledge of the structure of the model to be estimated. It is as if economic theory gave us direct access to the book of nature in which God had written down almost everything important, but somehow thought that it would be a good joke on people to leave out the values of the parameters. We do not have that sort of knowledge. We have to learn the structure of the model just as much as we have to learn the values of parameters from



empirical observation. Econometricians have frequently resisted the first use of statistical tests with a powerful, but ultimately vague and not-consistently-developed fear of data mining.

## V. Specification Search and Its Enemies

Among economists "data mining" is a pejorative term, nearly always invoked as a rebuke.

Unhappily, the metaphor has escaped them: gold mining is the *sine qua non* of uncovering treasure.

The economists' fear does have a basis. Imagine that we have a data-generating process such as

$$(7) \quad x_t = \delta + \varepsilon_t,$$

where  $\delta$  is a constant and the  $\varepsilon_t$  are i.i.d.  $N(0, \sigma^2)$ . Suppose that we seek to model this process with

### Model 3

$$(8) \quad x_t = \mu + \beta y_t + v_t,$$

where the  $v_t$  are i.i.d.  $N(0, \sigma^2)$  and  $y$  is some element of an infinite set of mutually independent, i.i.d. variables. Most elements of that set would prove to be insignificant as the regressor ( $y_t$ ) in (8) (i.e., we will not be able to reject the null hypothesis of  $\beta = 0$ ). But with a test size  $\alpha = 0.05$ , one time in twenty on average we will estimate a  $\hat{\beta}$  that rejects the null. If we follow a search procedure that allows us to keep searching until we find one of those cases, the probability of finding a significant regressor is one. This illustrates the *optional stopping problem* that is often thought to be the bane of hypothesis testing.

The optional stopping problem does not require that we have an infinite set of candidate variables. Even in a finite set the probability of finding significant regressors in a search procedure may be very far from the nominal size of the test used to evaluate their significance. In some cases, the probabilities can be calculated analytically. In more complex cases, they can be determined through simulations of the search procedure. To take one illustration, Lovell (1983, p. 4, Table 1)

considers a data-generating process like equation (7) and searches over a set of mutually orthogonal i.i.d candidate variables with a known variance for pairs in which at least one of the variables is significant in a model of the form

#### Model 4

$$(8) \quad x_t = \mu + \beta_1 y_{1t} + \beta_2 y_{2t} + v_t .$$

Table 1 shows that for a  $t$ -test with a size  $\alpha = 0.05$  the probability of the search procedure finding significant regressors (i.e., falsely rejecting the null implied in (7)) equals the test size only when there are only two candidate variables. As the number of candidate variables rises, the "true" significance level approaches unity. Lovell suggests that we penalize search by adapting critical values in line with the "true" significance levels rather than acting as if the nominal size of a single-shot test remained appropriate.

Two distinct costs accrue to not knowing the true parameters of the data-generating process (see Krolzig and Hendry 2001, p. 833; Hendry and Krolzig 2005, p. C40). The *cost of inference* is the uncertainty that arises from estimation in the case that we know the structure of the model. It is illustrated by the standard error of the estimates of  $\rho$  and  $\mu$  in Model 2. The *cost of search* is the cost that arises from the process of molding an econometric model into a form that accurately captures the salient features of the data-generating process. The take-home message of Lovell (1983) is that the costs of search are high, although in some cases calculable. The key lesson of Johansen's analysis of Models 1 and 2 is that the failure to mold the econometric model effectively may generate a large cost of inference: inferences based on Model 1 are systematically misleading about the likelihood of the mean of the data-generating process being close to 10. Another way to put this is that there is a *cost of misspecification* that offsets the cost of search and to evaluate any search procedure we have to adequately quantify the net costs.

In order to illustrate the failure of actual search procedures, Lovell (1983) conducts a more realistic simulation. He starts with a set of twenty actual macroeconomic variables. He then constructs nine models with different dynamic forms using subsets of the twenty as the independent variables in conjunction with definite parameter values and errors drawn from a random number generator. He then considers three search procedures over the set of twenty candidate variables: 1. stepwise regression; 2. maximizing  $\bar{R}^2$ ; and 3. max-min $|t|$  – i.e., choosing the set of regressors for which the smallest  $t$ -statistic in the set is the largest. Table 2 shows the empirically determined average type I and type II errors over fifty simulations of four of the models for a nominal test size of  $\alpha = 0.05$ . Since the relevant null hypotheses are that the coefficient on any variable is zero, type I error can be interpreted as falsely selecting a variable and type II error as falsely rejecting a variable. Each of the search procedures displays massive size distortions. The table also shows that type I and type II error are inversely related as intuition suggests.

It would violate logic to suggest that because these particular (and very simple) search procedures have poor properties that we should prefer not to search but simply to write down a model and to conduct a one-shot test. Illogical, yet one hears this procedure advised by colleagues from time to time. Johansen's example shows that the risks of misspecification vitiate that procedure. To his credit, Lovell does not suggest this, but instead suggests adjusting the nominal size of the tests to account for the degree of search. It also does not follow logically that, because these particular search procedures are poor, all search procedures are equally poor. The general prejudice against data mining captured in such phrases as "if you torture the data long enough, it will confess" are rather cavalier projections of the optional stopping problem in such simple cases as the one that Lovell examines to more complicated, but unanalyzed, situations. The problem with that analysis

and with the three simple search procedures in Table 2 is that the procedures themselves do not constitute a severe test of the specification.

An alternative approach to search is found in the so-called LSE approach of David Hendry and his colleagues. Hoover and Perez (1999) were the first to automate search procedures in this family. We showed, using an experimental design similar to Lovell's, that these procedures were in fact highly effective and not subject to the massive distortions that Lovell found with the three simple procedures (see also Hendry and Krolzig 1999). Hendry and Krolzig incorporated a refined version of Hoover and Perez's search procedure into a commercially available program, *PcGets*, where the name derives from one of its key characteristics that search is conducting from a general to a specific specification (Hendry and Krolzig 2005). Working with Hendry, Doornik developed a search algorithm in the same family that uses a substantially different approach to investigating the search paths (Doornik 2009). The algorithm, *Autometrics*, is now incorporated along with the econometrics package *PcGive* into the *Oxmetrics* econometrics suite.

Different in detail, all the procedures based on the LSE search methodology bear a strong family resemblance. Omitting many of the minor details, I will describe Hoover and Perez's (2003) search algorithm:

1. **Overlapping samples:** A search is conducted over two overlapping subsamples and only those variables that are selected in both subsamples are part of the final specification.
2. **General-to-specific simplification:** A general specification includes all the variables in the search universe as regressors. A subset of the variables (five in the results reported here) with the lowest  $t$ -statistics serve as starting points for simplification paths. To start on a path, one variable in this subset is deleted. The path is determined by a sequence of deletions, corresponding to the lowest  $t$ -statistic in the current specification until all the remaining variables are significant on test with size  $\alpha$ . At each deletion, the simplified regression is run through a battery of specification tests, including a subsample stability test and a test of the restrictions of the simplified model against the general model. If it fails a test, the variable is replaced and the variable with the next lowest  $t$ -statistic is deleted. The *terminal specification* is one in which either all variables are significant and the specification passes the battery of tests or in which no variable (significant or insignificant) can be removed without failing one of the tests in the battery.

3. **Selection among terminal specifications.** Tests are run among the terminal specifications to determine whether any one specification encompasses the others. If so, it is the *overall terminal specification* for the subsample. (see Mizon 1984; Mizon and Richard 1986 for a discussion of encompassing tests.) If not, a new specification is formed as the non-redundant union of the regressors of the terminal specifications, and the search procedure begins again along a single search path starting with this specification.
4. **Elimination of adventitious variables:** The *final specification* is the intersection of the regressors of the overall terminal specifications from the two subsamples.

Compared with the search algorithms investigated by Lovell, this is a complex procedure. Its general idea, however, is relatively simple. Just as Johansen's Model 2 nested Model 1, the initial general specification nests all possible final specifications. This guarantees that, if a model that adequately captures the data-generating process is nested in the general model, it will be possible to identify it in principle. Multiple search paths reduce the likelihood that low probability realizations will lead away from the target model. A criterion for the adequacy of the model is that it supports the statistical assumptions that would be maintained for purposes of inference, which include, for example, white noise errors, homoskedasticity, normality, and subsample stability (see Johansen 2006). The statistical tests in the search procedure measure how tightly these constraints are binding, and the algorithm uses the tests to mold the final specification, by eliminating possibilities that violate them.

The anti-data-mining rhetoric that is fueled by results such as those reported by Lovell would lead one to guess that such a test procedure would inevitably lead to wild distortions of size and power. But this is not a question in which it is wise to judge from the armchair. Hoover and Perez (2003) conducted a simulation study using a subset of the data used in Levine and Renelt's (1992) study of cross-country growth regressions: 36 variables  $\times$  107 countries. The dependent variable (an analogue to the *average rate of growth of GDP per capita 1960-89*, which was the target of their study) was constructed by selecting at random the independent variables. The coefficients for each variable were chosen by regressing *average rate of growth of GDP per*

*capita 1960-89* on the chosen independent variables. The simulation then creates an artificial dependent variable using error terms drawn from the residuals of this regression in the manner of a bootstrap. One hundred simulations were run for each of thirty specifications for true data-generating processes including for specifications involving 0, 3, 7, and 14 variables (12,000 specifications in all).

There is, of course, an irreducible cost of inference. Different simulations are parameterized with variables with wildly different signal-to-noise ratios. We know by construction that if our model were identical with the data-generating process, then the size of the test would be the same as the nominal size (assumed to be  $\alpha = 0.05$  in all the simulations). The *empirical size* is calculated as the ratio of the incorrect variables included to the total possible incorrect variables. The *size ratio* =  $\frac{\text{empirical size}}{\alpha}$  measures sins of commission. A size ratio of unity implies that search does not typically select variables that are not in the true model.

The power of the test depends on the signal-to-noise ratio. The *empirical power* for a given true variable is the fraction of the replications in which the variable is picked out by the search procedure; that is, it is the complement of the proportion of type II error. We determine the *true (simulated) power* through a bootstrap simulation of the data-generating process – that is, from the correct regressors without search. The true (simulated) power for a given true variable is the empirical power that one would estimate if there were no specification uncertainty, but sampling uncertainty remained. When the signal-to-noise ratio is low, the true (simulated) power will also be low; and, when it is high, the true (simulated) power will be high. The *power ratio* =  $\frac{\text{empirical power}}{\text{true (simulated) power}}$  measures sins of omission. A power ratio of unity

indicates that a search algorithm omits variables that appear in the true model only at the rate that they would fail to be significant if God whispered the true specification into your ear.

The two right-hand columns of Table 3 present the results for the general-to-specific search algorithm. The size ratios are very near to much below unity. Far from losing control over size in the manner of Lovell's various search algorithms, the general-to-specific procedure is *more* stringent than nominal size. Power ratios are close to unity. Given that size and power are inversely related, adjusting the nominal size of the underlying tests upward until they achieved a size ratio of unity would likely raise the power ratios towards unity as well.

The other four columns compare two other search algorithms that have been used in the literature on cross-country growth regressions and in other contexts. The left two columns refer to Leamer's (1983) extreme-bounds analysis as modified by Levine and Renelt (1992). Here each variable is taken in turn to be a *focus variable*. The focus variable is held fixed in regressions that include it and every possible three-variable subset of remaining variables. A 95-percent confidence interval is calculated for the focus variable for each of the regressions with different subsets of regressors. Any variable is eliminated as not robust if, any of these confidence intervals includes zero. The modified extreme-bounds analysis of Sala-i-Martin (1983) follows the same procedure, but treats a variable as not robust only if the confidence intervals include zero in more than 5 percent of the cases. Table 3 shows that the size ratios of the extreme-bounds analysis are tiny, implying that it almost never commits a sin of commission. But the power ratios are low and, in fact, fall to nearly zero when the number of regressors is large. In effect, its virtuous size is purchased with the wages of sins of omission: it simply rejects almost every regressor true or false. This is the problem that motivated the modified version evaluated in the middle two columns. Here the size ratios are very high, except when

there are no true variables to be found. This implies that the procedure is too permissive – many variables are selected that should not be. The power ratios are better behaved, though less well behaved than for the general-to-specific algorithm.

This simulation study shows that there are good and bad search procedures. A good search procedure is one in which the costs of search are low, so that all that remains are the costs of inference. The general-to-specific procedure appears to achieve that. This has been backed up by other simulation studies as well (see Hendry and Krolzig 2005; Doornik 2009). What accounts for the difference? I suggest that it is the severity of the testing procedure that arises from imposing multiple constraints on model through various specification tests. A theorem due to White (1990, pp. 379-380) clarifies the process. Informally, the theorem says: for a fixed set of specifications and a battery of specification tests, as the sample size grows toward infinity and increasingly smaller test sizes are employed, the test battery will – with a probability approaching unity – select the correct specification from the set. According to the theorem, both type I and type II error fall asymptotically to zero. Given sufficient data, only the true specification will survive a severe enough set of tests. The opponents of specification search worry that sequential testing will produce models that survive accidentally. Some hope to cure the problem through adjusting the critical values of statistical tests to reflect the likelihood of type I error. White's theorem, on the other hand, suggests that the true model is uniquely fitted to survive severe testing in the long run.<sup>6</sup> The key – as it is for breaking a code or solving a crossword puzzle – is to exploit the constraints of the data as fully as possible.

Asymptotic results are often suggestive but not determinative of what happens with fewer observations. The message, however, of the Monte Carlo simulations presented earlier is that it

---

<sup>6</sup> Recent analytical results for some specific aspects of search algorithms have added to our understanding and when and how they reduce the costs of search to a second-order problem; see Santos, Hendry and Johansen (2008) and Hendry and Johansen (2011).



is possible to design practical search algorithms that go a long way toward securing the promise of the asymptotic results. With models obtained through such severe search algorithms, the costs of search have been reduced sufficiently that it is reasonable to conduct inference as if we, in fact, knew the true model.

## References

- Boumans, Marcel. (2005) *How Economists Model the World in Numbers*. Abingdon, Oxford: Routledge.
- Doornik, Jurgen A. (2009) "Autometrics," in Jennifer L. Castle and Neil Shephard, editors. *The Methodology and Practice of Econometrics: A Festschrift in Honour of David F. Hendry*. Oxford: Oxford University Press, pp. 88-121.
- Giere, Ronald N. (1999) *Science Without Laws*. Chicago : University of Chicago Press, 1999.
- Giere, Ronald N. (2006) *Scientific Perspectivism*. Chicago: University of Chicago Press.
- Haavelmo, Trygve. (1944) "The Probability Approach in Econometrics," *Econometrica* 12(Supplement), iii-115.
- Hausman, Daniel M. (1992) *The Inexact and Separate Science of Economics*. Cambridge: Cambridge University Press.
- Heckman, James J. (2000) "Causal Parameters and Policy Analysis in Economics: A Twentieth Century Retrospective," *Quarterly Journal of Economics* 115(1), 45-97.
- Hendry, David F. (1987) 'Econometric methodology: a personal viewpoint', in Truman Bewley, editor. *Advances in Econometrics*, Vol. 2, Cambridge: Cambridge University Press.
- Hendry, David F. (2000) "Econometrics: Alchemy or Science?," in Hendry *Econometrics: Alchemy or Science*, 2<sup>nd</sup> edition. Oxford: Blackwell, pp. 1-28.
- Hendry, David F. and Hans-Martin Krolzig. (2005) "The Properties of Automatic GETS Modelling," *Economic Journal* 115(502) C32-C61.
- Hendry, David F., Søren Johansen (2011) "The Properties of Model Selection when Retaining Theory Variables," Discussion Paper No. 11-25, Department of Economics, University of Copenhagen.
- Hoover, Kevin D. (2006) "The Past as Future: The Marshallian Approach to Post-Walrasian Econometrics," in David Colander, editor. *Post Walrasian Macroeconomics: Beyond the Dynamic Stochastic General Equilibrium Model*. Cambridge: Cambridge University Press, 2006, pp. 239-257.
- Hoover, Kevin D. (2011) "Counterfactuals and Causal Structure," in Phyllis McKay Illari, Federica Russo, and Jon Williamson, editors. *Causality in the Sciences*. Oxford: Oxford University Press, pp. 338-360.

- Hoover, Kevin D. (forthcoming) "Pragmatism, Perspectival Realism, and Econometrics," in Aki Lehtinen, Jaakko Kuorikoski and Petri Ylikoski, editors. *Economics for Real: Uskali Mäki and the Place of Truth in Economics*. London: Routledge.
- Hoover, Kevin D. and Stephen J. Perez. (1999) "Data Mining Reconsidered: Encompassing and the General-to-Specific Approach to Specification Search,," *Econometrics Journal* 2(2), 1-25.
- Hoover, Kevin D. and Stephen J. Perez. (2003) "Truth and Robustness in Cross-country Growth Regressions," *Oxford Bulletin of Economics and Statistics* 66(5), 2004, pp. 765-798.
- Hoover, Kevin D. and Mark V. Sieglar. (2008) "Sound and Fury: McCloskey and Significance Testing in Economics," *Journal of Economic Methodology* 15(1), March 2008, 1-37
- Hoover, Kevin D. and Mark V. Sieglar. (2008) "The Rhetoric of 'Signifying Nothing': A Rejoinder to Ziliak and McCloskey," *Journal of Economic Methodology* 15(1), March 2008, 57-68.
- Hoover, Kevin D., Søren Johansen and Katarina Juselius. "Allowing the Data to Speak Freely: The Macroeconometrics of the Cointegrated Vector Autoregression," *American Economic Review* 98(2), 251-255.
- Johansen, Søren. 2006. "Confronting the Economic Model with the Data," in David Colander, editor. *Post Walrasian Macroeconomics: Beyond the Dynamic Stochastic General Equilibrium Model*. Cambridge: Cambridge University Press, pp. 287-300.
- Krolzig, H.-M. and Hendry, D. F. (2001). 'Computer Automation of General-to-Specific Model Selection Procedures', *Journal of Economic Dynamics and Control*, Vol. 25, pp. 831-866.
- Leamer, Edward E. (1983). "Let's Take the Con Out of Econometrics," *American Economic Review* 73(1), 31-43.
- Levine, Ross and David Renelt. (1992). "A Sensitivity Analysis of Cross-country Growth Regressions," *American Economic Review*. 82(4), 942-963.
- Liu, Ta-Chung. (1960) "Underidentification, Structural Estimation, and Forecasting," *Econometrica* 28(4), 855-865.
- Lovell, Michael C. (1983). Data mining. *Review of Economic Statistics* 65(1), 1-12.
- Mayo, Deborah. (1996) *Error and the Growth of Experimental Knowledge*. Chicago: University of Chicago Press.
- Mayo, Deborah G. and Aris Spanos. (2006) "Severe Testing as a Basic Concept in a Neyman-Pearson Philosophy of Induction," *British Journal for the Philosophy of Science* 57(2), 323-357.
- McCloskey, D.N. (2002) *The Secret Sins of Economics*, Chicago, IL: Prickly Paradigm Press.
- McCloskey, D.N., and Ziliak, S.T. (1996), "The Standard Error of Regressions," *Journal of Economic Literature*, 34(1), 97-114.
- Mizon, Grayham E. (1984) 'The Encompassing Approach in Econometrics', in David F. Hendry and Kenneth F. Wallis, editors. *Econometrics and Quantitative Economics*. Oxford: Basil Blackwell, pp. 135-72.

- Mizon, Grayham E. (1995) ‘Progressive Modelling of macroeconomic Time Series: The LSE Methodology’, in Kevin D. Hoover (ed.) *Macroeconometrics: Developments, Tensions and Prospects*, Boston: Kluwer, pp. 107–70.
- Mizon, Grayham E. and Jean-Francois Richard. (1986) “The Encompassing Principle and its Application to Testing Non-Nested Hypotheses,” *Econometrica* 54(3), 657-678.
- Morgan, Mary S. (2012) *The World in the Model: How Economists Work and Think*. Cambridge: Cambridge University Press.
- Morgan, Mary S. and Margaret Morrison. (1999) *Models as Mediators: Perspectives on Natural and Social Science*. Cambridge: Cambridge University Press.
- Santos, Carlos; David F. Hendry; and Søren Johansen. (2008) “Automatic Selection of Indicators in a Fully Saturated Regression,” *Computational Statistics* 23(2), 317-335,
- Sims, Christopher A. (1980) “Macroeconomics and Reality,” *Econometrica*, 48(1):1-48.
- Sterrett, Susan G. (2005) *Wittgenstein Flies a Kite: A Story of Models of Wings and Models of the World*. New York: Pi Press.
- White, Halbert (1990) “A Consistent Model Selection Procedure Based on  $m$ -testing,” in C.W.J. Granger, editors. *Modelling Economic Series: Readings in Econometric Methodology*. Oxford: Clarendon Press, 369–383.
- Ziliak, S.T., and McCloskey, D.N. (2004), “Size Matters: The Standard Error of Regressions in the American Economic Review,” *Journal of Socio-Economics*, 33(5), 527–546.

**Table 1**  
**The Dependence of the True Size of a Hypothesis Test on Search**

<b>Number of Variables in Pool</b>	<b>True Significance Level</b>
2	0.050
5	0.120
10	0.226
20	0.401
100	0.923
500	0.999

Notes: Table is based on Lovell (1983, Table 1). Variables in the pool are independent i.i.d and the hypothesis of no relationships with the dependent variables is true. Search procedures regresses dependent variable on pairs of variables in the pool until at least one coefficient is statistically significant at the  $\alpha = 0.05$  level. The *true significance level* is the actual proportion of searches in which a significant regressor is identified.

**Table 2**  
**Error Rates for Three Simple Search Algorithms**

<b>Error Rates (percent)</b>			
	<b>Stepwise Regression</b>	<b>max <math>\bar{R}^2</math></b>	<b>max-min <math> t </math></b>
<b>Type I Error</b>	30	53	81
<b>Type II Error</b>	15	8	0

Notes: Based on Lovell (1983, Table 7). The table reports the average error rates over 50 simulations of four models using three search algorithms.

**Table 3. The Efficacy of Three Search Algorithms**

<b>Models with:</b>	<b>Extreme-bounds analysis</b>		<b>Modified Extreme-bounds analysis</b>		<b>General-to-Specific</b>	
	<b>Size Ratio<sup>1</sup></b>	<b>Power Ratio<sup>2</sup></b>	<b>Size Ratio<sup>1</sup></b>	<b>Power Ratio<sup>2</sup></b>	<b>Size Ratio<sup>1</sup></b>	<b>Power Ratio<sup>2</sup></b>
<b>0 true variables</b>	0.060		1.10		0.75	
<b>3 true variables</b>	0.003	0.43	5.17	0.77	0.77	0.95
<b>7 true variables</b>	0.030	0.13	5.89	1.10	0.81	0.93
<b>14 true variables</b>	0.020	0.04	5.45	0.67	1.02	0.82

Notes: The table is original Table 1 in Hoover and Perez (2003). The basic data are a pool of 36 variables described in Memorandum 1 downloadable from our websites (<http://www.econ.ucdavis.edu/faculty/kdhoover/research.html> and <http://www.csus.edu/indiv/p/perezs/Data/data.htm>). For each number of true variables, 30 models are specified by choosing the indicated number of regressors at random from the pool. Coefficients are calibrated from a regression of the chosen regressors on the actual average growth rate. 100 dependent variables are created from the same regressors and coefficients and error terms constructed with a wild bootstrap procedure from the errors of the calibrating regression. Specification searches are then conducted by each of the three methods and the number of type I and type II errors are recorded. Statistics reported here average over each of the 100 simulations for each of the 30 models. Details of the simulations and the search procedures are found in Section 2 and Appendix 1.

<sup>1</sup>Size is calculated as the proportion of incorrect variables included (significantly for general-to-specific) to the total possible incorrect variables. The Size Ratio is average ratio of the size to the nominal size (0.05) used as the critical value in all the hypothesis tests in the search procedures. A Size Ratio of 1.00 indicates that on average the size is equal to the nominal size (0.05).

<sup>2</sup>Power is calculated as the proportion of times a true variables is included (significantly for the general-to-specific procedure). The true (simulated) power is based on the number of type II errors made in 100 simulations of the true model without any search. The Power Ratio is the average ratio of power to true (simulated) power. A Power Ratio of 1.00 indicates that on average the power is equal to the true (simulated) power. The Power Ratio is not relevant when there are no true variable