

# From Association to Causation via Regression\*

David Freedman

*Statistics Department, University of California, Berkeley, California 94720*

Received October 3, 1995; revised September 1996

For nearly a century, investigators in the social sciences have used regression models to deduce cause-and-effect relationships from patterns of association. Path models and automated search procedures are more recent developments. In my view, this enterprise has not been successful. The models tend to neglect the difficulties in establishing causal relations, and the mathematical complexities tend to obscure rather than clarify the assumptions on which the analysis is based. Formal statistical inference is, by its nature, conditional. If maintained hypotheses  $A, B, C, \dots$  hold, then  $H$  can be tested against the data. However, if  $A, B, C, \dots$  remain in doubt, so must inferences about  $H$ . Careful scrutiny of maintained hypotheses should therefore be a critical part of empirical work—a principle honored more often in the breach than the observance. This paper focuses on modeling techniques that seem to convert association into causation. The object is to clarify the differences among the various uses of regression, as well as the source of the difficulty in making causal inferences by modeling. The discussion will proceed mainly by examples, ranging from Yule (*J. R. Stat. Soc.* **62** (1899), 249–295) to Spirtes, Glymour, and Scheines (“Causation,” *Lect. Notes in Statist.* Vol. 81, Springer-Verlag, New York/Berlin, 1993). © 1997 Academic Press

## 1. OUTLINE

Many treatments of regression seem to take for granted that the investigator knows the relevant variables, their causal order, and the functional form of the relationships among them; measurements of the independent variables are assumed to be without error. Indeed, Gauss developed and used regression in physical science contexts where these conditions hold, at least to a very good approximation.<sup>1</sup> Today, the textbook theorems that justify regression are proved on the basis of such assumptions.

\* Presented at the Notre Dame Conference on Causality in Crisis, Oct. 15–17, 1993.

<sup>1</sup> Gauss was fitting orbits to astronomical observations, with least squares to estimate the elements of the orbits [21]. Stigler [64, pp. 145–146] awards priority to Legendre [36].

In the social sciences, the situation seems quite different. Regression is used to discover relationships or to disentangle cause and effect. However, investigators have only vague ideas as to the relevant variables and their causal order; functional forms are chosen on the basis of convenience or familiarity; serious problems of measurement are often encountered.

Regression may offer useful ways of summarizing the data and making predictions. Investigators may be able to use summaries and predictions to draw substantive conclusions. However, I see no cases in which regression equations, let alone the more complex methods, have succeeded as engines for discovering causal relationships. Of course, there may be success stories that I have not found; nor does a track record of failure necessarily project into the future.

One of the first applications of regression techniques to social science is Yule [71]. Recent examples will be found in Spirtes, Glymour, and Scheines [62], to be cited here as SGS. (The SGS theory is summarized in Glymour [23], cited as CG.) SGS have attracted considerable attention in the philosophy of science, because they have developed computerized algorithms that search for path models. With their algorithms, SGS claim to make rigorous inferences of causation from association. This is a bold claim, which does not survive examination.

The balance of this paper is organized as follows. Section 2 discusses Yule's work. Sections 3 and 4 explain the critical data of "exogeneity." Section 5 describes a contemporary regression model. Sections 6–10 review SGS and reanalyze some of their examples. Sections 11–12 canvass some mathematical issues. Possible responses to my critique will be found in Section 13. There is a brief review of the literature in Section 14, and conclusions are presented in Section 15. For ease of reference, standard formulas for regression are given in an appendix. I have tried to make most of the paper accessible to nonstatistical readers, particularly if they will permit the occasional undefined technical term; Sections 11 and 12 are more specialized.

## 2. YULE'S REGRESSION MODEL FOR PAUPERISM

One of the first regression models in social science was developed by Yule—"An Investigation into the Causes of Changes in Pauperism in England, Chiefly During the Last Two Intercensal Decades."<sup>2</sup> In late 19th century England, poor people could be supported either inside the poor house or outside. Did provision of support outside the poor house increase the number of poor people?

<sup>2</sup> See [71; 64, pp. 345–358; 11].

To address this issue, Yule used data from the censuses of 1871, 1881, and 1891. (In England, the census is taken in years that end with 1.) He considered the periods 1871–1881 and 1881–1891, relating changes in the number of paupers to changes in the “outrelief ratio,” that is, the ratio between the number of paupers supported outside the poor house and inside. He used regression to control for two confounders—changes in the population and its age structure.

His equation can be written as follows:

$$\Delta Paup = a + b \times \Delta Out + c \times \Delta Pop + d \times \Delta Old + error. \quad (1)$$

Here,  $\Delta$  stands for percentage difference, *Paup* for the number of paupers, *Out* for the outrelief ratio, *Pop* for population size, and *Old* for the proportion of people aged 65 and over.

Yule’s unit of analysis was the “union,” which seems to have been a small geographical area like a county.<sup>3</sup> He had four kinds of areas: rural, mixed, urban, metropolitan. He used “Ordinary Least Squares” (OLS) to estimate the coefficients from the data, with a “50 cm. Gravet” slide rule to do the arithmetic.

To be more specific, Yule estimated a separate equation for each time period (1871–81 and 1881–91) and each kind of area. There were 2 time periods and 4 kinds of areas, thus,  $2 \times 4 = 8$  equations. Within a time period, all areas of the same kind—for instance, all rural unions—are governed by one equation. (By coincidence, there are 4 coefficients in each equation, and 4 kinds of areas.)

Yule was looking for the “Hooke’s Law of Poverty.” Nature ran an experiment, with lots of variation over time and geography, and Yule analyzed the results. Regression was needed to control for the confounding effects of change in population and age structure. The equations were held to show that, other things being equal, changes in the outrelief ratio create corresponding changes in the number of paupers. Indeed, if you increase the outrelief ratio by one percentage point but hold the other factors constant, you will increase the number of paupers by  $b$  percent,  $b$  being the coefficient of  $\Delta Out$  in Eq. (1). More qualitatively, if  $b$  is positive, welfare creates paupers.

For a moment, I turn from Yule to methodology. A regression equation like (1) is usually written as

$$Y = X\beta + \varepsilon. \quad (2)$$

In this equation, the vector  $Y$  represents the dependent variable, like pauperism; the matrix  $X$  represents the explanatory (or “independent”)

<sup>3</sup> There were about 600 such areas in England. A poor-law union “consisted of two or more parishes combined for administrative purposes.” [64, p. 346].

variables, like the outrelief ratio, population, and age structure. These are observable. The vector  $\beta$  represents parameters, which are not observable but may be estimated from the data; parameters are “social constants,” which characterize the process that generated the data. In Yule’s equation,  $\beta$  has four components—the parameters  $a, b, c, d$  in Eq. (1). The error or “disturbance” term  $\varepsilon$  is also unobservable and represents the impact of chance factors unrelated to  $X$ . Statistical inferences are often based on “stochastic assumptions” about  $\varepsilon$ ; e.g.,  $\varepsilon$  is independent of  $X$  and its components are independent and identically distributed with mean 0. For details, see the Appendix.

Three possible uses for regression equations are

- (i) to summarize data, or
- (ii) to predict values of the dependent variable, or
- (iii) to predict the results of interventions.

Yule could certainly have summarized his data by saying that for a given time period and unions of a specific type, with certain values of the explanatory variables, the change in pauperism was about so much and so much. In other words, he could have used his equations to estimate the average value of  $Y$ , given the values of  $X$ . This use of regression may run into technical problems if there are outliers, or nonlinearities in the regression surface. However, at least in principle, there do seem to be technical fixes for such problems. Furthermore, stochastic assumptions about the disturbance term play almost no role. Therefore, like most statisticians, I believe that regression can be quite helpful in summarizing large data sets.

For prediction, there is a *ceteris paribus* assumption: the system will remain stable. Prediction is already more complicated than description. On the other hand, if you make a series of predictions and test them against data, it may be possible to show that the system is stable, or sufficiently stable for regression to be quite helpful.<sup>4</sup> Again, any particular use of regression to make predictions may go off the rails, but there do not seem to be essential difficulties of principle involved.

Causal inference is different, because a change in the system is contemplated; for example, there will be an intervention. Descriptive statistics tell you about the correlations that happen to hold in the data; causal models claim to tell you what will happen to  $Y$  if you change  $X$ . Indeed, regression is often used to make counterfactual inferences about the past: what would  $Y$  have been if  $X$  had been different? This use of regression

<sup>4</sup> Meehl [41] provides some well-known examples. Predictive validity is best demonstrated by making real *ex ante* forecasts in several different contexts: see Ehrenberg and Bound [13].

to make causal inferences is the most intriguing—and the most problematic. Difficulties are created by omitted variables, incorrect functional form, etc. Of course, if the results of causal modeling were with any frequency checked against the results of interventions, the balance of argument might be very different.<sup>5</sup>

For description and prediction, the numerical values of the individual coefficients fade into the background; it is the whole linear combination on the right-hand side of the equation that matters. For causal inference, it is the individual coefficients that do the trick. In Eq. (1), for example, it is  $b$  that should tell you what happens to pauperism when the outrelief ratio is manipulated.

At this remove, the flaws in Yule's argument may be apparent. For example, there seem to be some important variables missing from the equation, including variables that measure economic activity. Here is Yule's comment on the last-named factor [71, p. 253]:

A good deal of time and labour was spent in making trial of this idea, but the results proved unsatisfactory, and finally the measure was abandoned altogether.

Yule [71] seems to have used the rate of population growth— $\Delta Pop$  in Eq. (1)—as a proxy for economic activity, although that creates ambiguity. Other things being equal, population growth will by itself add to the number of paupers; in its role as proxy, however, population growth should reduce pauperism.

The equations for metropolitan unions are shown below, for 1871–1881 and 1881–1891.<sup>6</sup>

(1871–1881)

$$\Delta Paup = 13.19 + 0.755 \times \Delta Out - 0.322 \times \Delta Pop \\ - 0.022 \times \Delta Old + \text{residual}.$$

(1881–1891)

$$\Delta Paup = 1.36 + 0.324 \times \Delta Out - 0.369 \times \Delta Pop \\ + 1.37 \times \Delta Old + \text{residual}.$$

For example, one metropolitan union is Westminster. Over the period 1871–1881, the percentage changes in *Out*, *Pop*, and *Old* are  $-73$ ,  $-9$ ,

<sup>5</sup> Also see Manski [40].

<sup>6</sup> These, and the other six equations, are reported in Yule [71, Table C, p. 259]. His Table XIX gives data for metropolitan unions, in the form of “percentage ratios” for 1871–1881 rather than differences, apparently to avoid negative numbers. The equations were fitted to data; the numerical coefficients in the displays are estimates for the corresponding parameters in (1); the residuals are observable, but are only approximations to unobservable disturbance terms.

and 5, respectively. The percentage change in *Paup* predicted from the regression equation is

$$13.19 + 0.755 \times (-73) - 0.322 \times (-9) - 0.022 \times 5 = -39.$$

The actual percentage change in *Paup* is  $-48$ . The “residual” is

$$\text{residual} = \text{actual} - \text{predicted} = -48 - (-39) = -9.$$

The coefficients in the regression equation are estimated so as to minimize the size of the residuals. (Technically, it is the sum of the squares that is minimized—hence the term “least squares.”) The linear combination of explanatory variables on the right side of the equation has therefore been optimized; but there is no guarantee that individual coefficients will make much sense.

There are some noticeable inconsistencies in Yule’s coefficients, over time and across the various kinds of geography. Nor are the signs of the coefficients entirely reasonable. These inconsistencies may not by themselves be fatal, but they certainly raise the question of whether the equations hold true for any well-defined population of times and places. If the coefficients do not have a life of their own—outside Yule’s particular data set—they cannot be used to answer questions of the form, “What would happen if you change the outrelief ratio?” The coefficients may be useful for descriptive purposes, but not for causal inference or even prediction.

Moreover, there are familiar difficulties of interpretation. At best, Yule showed that changes in pauperism and the outrelief ratio were associated, even after adjusting for changes in the population and its age structure. The direction of the causal arrow, however, is by no means clear. Yule’s theory is that outrelief is the cause and pauperism is the effect. That is a reasonable view. However, the opposite idea seems equally tenable—a union that is flooded with paupers may not be able to build poor houses fast enough and resorts to outrelief. If so, pauperism causes outrelief. Also, Governor Pete Wilson’s theory may have some plausibility for 19th century England if not 20th century California: unions that provide generous outrelief attract paupers from elsewhere.<sup>7</sup>

Yule must have been aware of these problems. After allocating the changes in pauperism to their various causes (including the residual), he

<sup>7</sup> According to Stigler [64, pp. 356–357], Pigou criticized Yule for ignoring “the non-quantitative facts of the situation . . . . It is well known that, during recent years, those unions in which out-relief has been restricted have, on the whole, enjoyed a general administration much superior to that of other unions.” Stigler responds that “Pigou’s ad hoc speculation . . . could not, of course, be disproved from the data Yule used.” In effect, this allows Yule to defend himself by pleading ignorance.

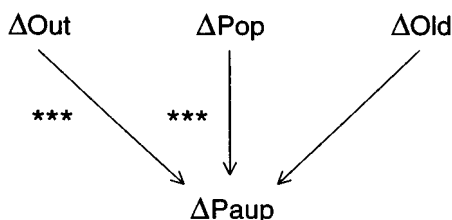


FIG. 1. Yule's model for pauperism. The figure represents Eq. (1) in graphical form. The asterisks denote a high degree of statistical significance. To determine the asterisks, I recomputed Yule's regression for the metropolitan unions over the period 1871–1881, using data in his Table XIX. I replicated his coefficients, as shown in the display, although roundoff error is quite large:

$$\Delta \text{Paup} = 12.884 + 0.752 \times \Delta \text{Out} - 0.311 \times \Delta \text{Pop} + 0.056 \times \Delta \text{Old} + \text{residual.}$$

10.367	0.135	0.067	0.223
1.24	5.57	– 4.645	0.25

Under the coefficients are standard errors (SEs) and *t*-statistics. The SE indicates the likely size of the difference between an estimated coefficient and its true value. The *t*-statistic is the ratio of an estimate to its SE. Generally, a *t*-statistic above 2 or 3 in absolute value indicates that the corresponding parameter is unlikely to be truly 0. The parameters are features of the model, and the SEs are computed on the basis of the stochastic assumptions in the model; for details, see the appendix. In Fig. 1, the explanatory variables are correlated; such correlations are often signaled by curved, double-headed arrows; error terms are not shown either.

withdraws all causal claims with one deft sentence: “Strictly, for ‘due to’ read ‘associated with.’” [71, p. 270, footnote 25]. Yule's paper is quite modern in spirit, with two exceptions: he did not rely on statistical significance, and he did not use a graph. Figure 1 brings him up to date.

### 3. REGRESSION ESTIMATES AND CONDITIONAL EXPECTATIONS

In the regression model (2), *Y* is the dependent variable, like pauperism; *X* represents the explanatory variables, like the outrelief ratio, population, and age structure. If all goes well, the regression equation will estimate the “conditional expectation” of *Y* given *X* = *x*, that is, the average value of *Y* corresponding to given values for the explanatory variables.

To clarify the definitions, consider two procedures:

*Procedure 1.* Select subjects with *X* = *x*; look at the average of their *Y*'s.

*Procedure 2.* Intervene and set  $X = x$  for some subjects; look at the average of their  $Y$ 's.

These procedures are quite different. The first involves the data set as you find it. The second involves an intervention.

Regression does seem to let you move from selection to intervention; that is why the technique is so popular. However, regression approximates the selection procedure, rather than intervention. Nor does the statistical analysis prove that the two procedures give the same results; how could it? Instead, causal inferences are made by *assuming* that selection tells you what would happen if you were to intervene.

The phrase “ $X$  is exogenous” is often taken to mean that selecting on  $X$  will produce the same results as intervening to set the value of  $X$ —the basic assumption in many analyses. Exogeneity also has weaker meanings, to be taken up later. The ambiguity is unfortunate, because analysts may assume exogeneity in a weak sense and proceed as if they had established something more. It is only exogeneity in the strong sense defined above that enables you to predict the results of interventions from nonexperimental data.

The distinction between selection and intervention is acknowledged even by the modelers (Pearl [44, p. 396]):

Formally speaking, probabilistic analysis is indeed sensitive only to covariations, so it can never distinguish genuine causal dependencies from spurious correlations . . . .

Such admissions—like Yule’s [71] footnote 25—are fatal to the enterprise. Of course, Pearl does not give up. For instance, he goes on to say that experiments just provide the opportunity to observe yet more correlations, a move he attributes to Simon [59].

Figure 2 is Pearl’s [44]. On the left, it seems that  $X$  and  $Z$  cause  $Y$ ; manipulating  $X$  or  $Z$  will change  $Y$ . However, if only we had measured the

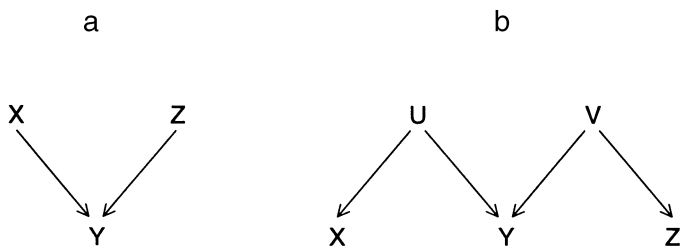


FIG. 2. After Judea Pearl [44, p. 397]. Causation cannot be inferred from association by using causal models. In panel (a),  $X$  and  $Z$  are assumed to be independent. In panel (b),  $U$  and  $V$  are assumed to be independent; it may be shown in consequence that  $X$  and  $Z$  are independent. Also see Duncan [12, pp. 113–127].



variables  $U$  and  $V$ , we might have seen that they were the joint causes of  $X$ ,  $Y$ , and  $Z$ , as in the right-hand panel. If so, manipulating  $X$  and  $Z$  will not change  $Y$  at all. No amount of statistical analysis on the observables—on  $X$ ,  $Y$ , and  $Z$ —can tell us which panel expresses the right theory. Indeed, matters can be arranged so that both theories lead to the same joint distribution for the observables.

#### 4. TWO IDEAS OF CONDITIONAL PROBABILITIES

The distinction between the two ideas of conditioning—selecting subjects with  $X = x$ , or intervening to set  $X = x$ —seems fundamental. A concrete example may help, and conditional probabilities are easier to deal with than conditional expectations.

Many studies have demonstrated an association between cervical cancer and exposure to two sexually transmitted diseases—herpes and chlamydia. Suppose we had data as shown in Table I. The incidence rate of cervical cancer is 200 per 100,000 for women exposed to herpes and chlamydia (top left); 116 per 100,000 for women exposed to herpes but not chlamydia; and 130 per 100,000 for those exposed to herpes, the two exposure categories for chlamydia being combined. Other cells may be read in a similar way.

With sample data, there is a role for technical statistics in estimation and testing—for instance, to see if the rates within a row are constant across columns. However, the real question is not association but causation. Does herpes cause cervical cancer? What about chlamydia? Biotechnology might find a way to eliminate *Herpes simplex* as well as *Chlamydia trachomatis*. That would be a great relief, but would it reduce the incidence rate of cervical cancer?

To consider the issue of causality more directly, suppose that we actually know the rates for the population of interest, as shown in Table I. Statistical testing must now fade into the background. The overall inci-

TABLE I  
Rate of Cervical Cancer Cases per 100,000 Women, by Exposure to Chlamydia and Herpes

	Chlamydia		Total
	Yes	No	
Herpes			
Yes	200	116	130
No	180	80	87
Total	190	90	100

Note. Data are hypothetical.

dence rate is 100 cervical cancers per 100,000 women (Table I, bottom right). Among women exposed neither to herpes nor to chlamydia, the rate is lower—80 per 100,000. If cervical cancer is caused by herpes and chlamydia, eliminating the microorganisms responsible for those diseases should reduce the incidence rate of cervical cancer from 100 to 80 per 100,000. On the other hand, if the relationship is not causal, eliminating those microorganisms will have little effect on the incidence rate of the cancer.

To be more explicit, 80/100,000 has been found by selecting women who are exposed to neither herpes nor chlamydia and by computing the incidence rate of cervical cancer for that group, one interpretation of conditional probability. If we intervene and eliminate the two diseases, we want to know the rate after the intervention; that is another interpretation. The two interpretations are different, because the underlying procedures are different. Statistical analysis of the numbers in the table, however refined or complex, cannot prove that a hypothetical intervention will give the same results as selection. This may seem obvious, even banal; but if you grant the point, the causal modeling game is largely over.

What is the situation for Table I? The story is far from certain. Current epidemiological opinion favors the idea that cervical cancer is caused by certain strains of human papilloma virus (HPV); herpes and chlamydia have no etiologic role, but serve only as markers for exposure to HPV. If that opinion is correct, wiping out herpes and chlamydia will have no impact on rates of cervical cancer.

Due in part to the rarity of cervical cancer, cohort studies do not seem to be available. (The numbers in Table I, although hypothetical, are not unreasonable.) My point is even stronger for the real studies of the association between cervical cancer and herpes or chlamydia. Problems created by incomplete data cannot simplify the task of inferring causation from association.<sup>8</sup>

## 5. ANOTHER REGRESSION EXAMPLE

Rindfuss *et al.* [55] propose a model to explain the process by which a woman decides how much education to get, and when to have her first child. The model illustrates many features of contemporary technique.<sup>9</sup>

<sup>8</sup> For a discussion of the epidemiology, see Cairns [4], Peto and zur Hausen [51], Sherman *et al.* [58], Hakama *et al.* [25], Muñoz *et al.* [75].

<sup>9</sup> I use this example because it is discussed by SGS [62, pp. 139–140].

Before we take up the model, let the authors say what they were trying to do:

The interplay between education and fertility has a significant influence on the roles women occupy, when in their life cycle they occupy these roles, and the length of time spent in these roles. . . . This paper explores the theoretical linkages between education and fertility. . . . It is found that the reciprocal relationship between education and age at first birth is dominated by the effect from education to age at first birth with only a trivial effect in the other direction. [Abstract]

No factor has a greater impact on the roles women occupy than maternity. Whether a woman becomes a mother, the age at which she does so, and the timing and number of subsequent births set the conditions under which other roles are assumed. . . . Education is another prime factor conditioning female roles. [p. 431, footnote omitted]

The overall relationship between education and fertility has its roots at some unspecified point in adolescence, or perhaps even earlier. At this point aspirations for educational attainment as a goal in itself and for adult roles that have implications for educational attainment first emerge. The desire for education as a measure of status and ability in academic work may encourage women to select occupational goals that require a high level of educational attainment. Conversely, particular occupational or role aspirations may set standards of education that must be achieved. The obverse is true for those with either low educational or occupational goals. Also, occupational and educational aspirations are affected by a number of prior factors, such as mother's education, father's education, family income, intellectual ability, prior educational experience, race, and number of siblings. [p. 432, citations omitted]

The model used by Rindfuss *et al.* [55] is shown in Fig. 3. The diagram corresponds to two linear equations in two unknowns, ED and AGE (variables are defined in Table II):

$$ED = a \times AGE + A, \quad (3)$$

$$AGE = a' \times ED + A'. \quad (4)$$

According to the model, a woman chooses her educational level and age at first birth as if by solving these two equations for the two unknowns.

The coefficients  $a$  and  $a'$  are "social constants," to be estimated from the data. The terms  $A$  and  $A'$  take background factors into account:

$$A = A_0 + b \times DADSOCC + c_1 \times RACE + \cdots + c_7 \times YCIG \\ + \text{random error drawn from a box}, \quad (5)$$

$$A' = A'_0 + b' \times FEC + c'_1 \times RACE + \cdots + c'_7 \times YCIG \\ + \text{another random error drawn from a box}. \quad (6)$$

Again, the parameters  $A_0, b, c_1, \dots$  are social constants to be estimated from the data. The random errors are assumed to have mean 0, to be

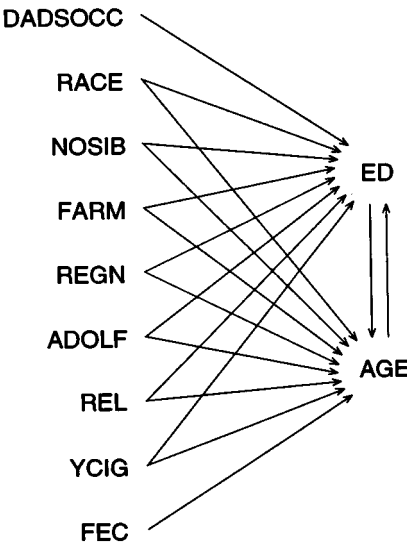


FIG. 3. The model in diagram form [55; 62, p. 140]. Variables are defined in Table II. Explanatory variables (DADSOCC, RACE, etc.) are correlated; error terms are not shown in the diagram.

TABLE II  
Variables in the Model [55]

	<i>The endogenous variables</i>
ED	Respondent's education (Years of schooling completed at first marriage)
AGE	Respondent's age at first birth
	<i>The exogenous variables</i>
DADSOCC	Respondent's father's occupation
RACE	Race of respondent (Black = 1, other = 0)
NOSIB	Respondent's number of siblings
FARM	Farm background (coded 1 if respondent grew up on a farm, else 0)
REGN	Region where respondent grew up (South = 1, other = 0)
ADOLF	Broken family (coded 0 if both parents present at age 14, else 1)
REL	Religion (Catholic = 1, other = 0)
YCIG	Smoking (coded 1 if respondent smoked before age 16, else coded 0)
FEC	Fecundability (coded 1 if respondent had a miscarriage before first birth; else coded 0)

*Note.* The data are from a probability sample of 1766 women 35–44 years of age residing in the continental United States; the sample was restricted to ever-married women with at least one child. DADSOCC was measured on Duncan's scale, combining information on education and income; missing values were imputed at the overall mean. SGS [62, p. 139] gives the wrong definitions for NOSIB and ADOLF.

statistically independent from woman to woman, and to be identically distributed. Correlations across Eqs. (5) and (6) are permitted.

Equations (3)–(6) are not quite regression equations, due to the simultaneity of (3) and (4); fitting by OLS (ordinary least squares) would create “simultaneity bias.” Thus, Rindfuss *et al.* [55] use an estimation procedure called “two-stage least squares.”<sup>10</sup> FEC does not enter into Eq. (5), nor DADSOCC into Eq. (6). Graphically, there is no arrow from DADSOCC to AGE in Fig. 3; likewise, there is no arrow from FEC to ED. These behavioral assumptions are critical to the statistical enterprise. Without them, or some similar assumptions, two-stage least squares could not be used. Technically, the system would not be “identifiable” (Section 11.4).

The main empirical finding: The estimated coefficient of AGE in the first equation is not “statistically significant”; i.e., the coefficient  $a$  in (3) could be zero. The sort of woman who drops out of school to have a child would drop out anyway.

If looked at coldly, the argument may seem implausible. A critique can be given along the following lines:

(i) *Statistical assumptions.* Just why are the errors independent and identically distributed across the women? Independence may be reasonable, but heterogeneity is more plausible than homogeneity.

(ii) *The assumption of constant coefficients.* Rindfuss *et al.* [55] are assuming that the same parameters apply to all women alike, from poor blacks in the cities of the Northeast to rich whites in the suburbs of the West. Why?

(iii) *Omitted variables.* Surely, important variables have been omitted from the model, including two that were identified by Rindfuss *et al.* [55]—aspirations and ability. Malthus thought that wealth was an important factor. Social class matters, and DADSOCC measures only one of its aspects.<sup>11</sup>

(iv) *What about the “no arrow” assumptions, from DADSOCC to AGE and FEC to ED?*

(v) *Are FEC and DADSOCC exogenous?*

(vi) *Are the equations “structural”?*

Questions (iv)–(vi) will be discussed in the next section, as will the idea of “structural” equations.

<sup>10</sup> See, e.g., Maddala [39]; for discussion, see Daggett and Freedman [9].

<sup>11</sup> The solution to the “omitted variable” problem may seem easy—just throw some more variables into the model. The difficulties are explored in Clogg and Haritou [6]. Also see Freedman [17].

### 5.1. *A Thought Experiment*

A simpler version of the model restricts attention to a more homogeneous group of women, where the only relevant background factors are DADSOCC and FEC. To make causal inferences from the data using the model, we need to believe that the arrows are as shown in Fig. 4, that DADSOCC and FEC are exogenous, and that the equations are “structural.” The following thought experiment may help to define the last term, and the empirical commitments behind the words.<sup>12</sup>

The *gedanken* experiment involves two groups of women. In both groups, fathers are randomized to jobs, and some of the daughters are chosen at random to have a miscarriage before their first child. (The statistical terminology of randomization is dry; the *gedanken* experimentalist intervenes, for instance, to make the fathers do one job rather than another: professors are caused to work as plumbers, and taxi drivers are installed as hospital anesthetists.)

*Group I.* Daughters are randomized to the various levels of ED, and AGE is observed as the response. (The *gedanken* experimentalist strikes again, forcing some women to stay in school longer than they wish, while preventing others from continuing their education.)

*Group II.* Daughters are randomized to the various levels of AGE, and ED is observed as the response. (More *gedanken* intervention is needed.)

The statistical model can now be translated. For the women in Group I, AGE should not depend on DADSOCC—the “no arrow” assumption; however, AGE should depend linearly on ED. For the women in Group II, ED should not depend on FEC—the other “no arrow” assumption;

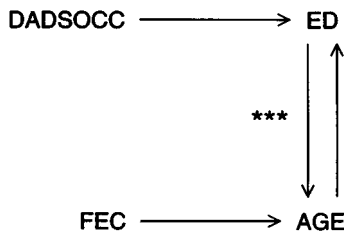


FIG. 4. A simpler version of the model.

<sup>12</sup> Also see Pearl [46, 47].

however, ED should depend linearly on DADSOCC. Rindfuss *et al.*'s [55] discovery is that ED would not depend on AGE.

There is one final assumption: the equations and parameters that describe the responses of the women in the experiment must also describe the natural situation. That is what "structural" means. For instance, a woman who freely chooses her educational level and her time to bear children does so by using the same equations as a woman made to give birth at a certain age. In short, with respect to the matters at issue, life in Des Moines proceeds more or less along the same lines as life in the Gulag.

The thought experiment provides the intellectual foundation for the model, by articulating the background assumptions. These assumptions have not been subjected—cannot be subjected—to direct empirical proof, nor can assumptions be validated by appealing to thought experiments that are almost unthinkable. Do the modelers have some other method in reserve? If the assumptions remain unvalidated, what is the logical status of their implications?

## 5.2. *Exogeneity*

Identifying the exogenous variables is a major problem. For example, you can obtain results quite different from those of Rindfuss *et al.*, [56] by using variables other than DADSOCC and FEC as "instruments."<sup>13</sup> Rindfuss *et al.* [56, pp. 981–982] respond that estimates made by

instrumental variables... require strong theoretical assumptions... and can give quite different results when alternative assumptions are made... it is usually difficult to argue that behavioral variables are truly exogenous and that they affect only one of the endogenous variables but not the other.

In short, results can depend quite strongly on assumptions of exogeneity, and there is no good way to justify one set of assumptions rather than another. Also see Bartels [1], who comments on the impact of exogeneity assumptions and the difficulty of verification.

<sup>13</sup> See Hofferth and Moore [27, 42]. An "instrument" is an exogenous variable, used as part of the two-stage least squares estimation procedure. Some investigators may draw a terminological distinction: an "instrument" is exogenous, but does not appear as an explanatory variable in the equation being estimated. For purposes of estimation, exogenous variables are assumed to be independent of error terms; this does not suffice for causal inference (Section 11). Even the independence assumption is not to be made lightly: see Clogg and Haritou [6].

## 6. AUTOMATED SEARCHES FOR CAUSALITY

SGS [62] have computerized algorithms that search for path models. Using the algorithms, SGS claim to make rigorous inferences of causation from association. For present purposes, a “path model” is a recursive system of regression equations, in which the dependent variables from some equations are used as explanatory variables in later equations.<sup>14</sup>

The basic idea in path models is this: putative causes combine with parameters and random errors by multiplication and addition in order to produce their effects. I have discussed such models elsewhere and do not believe they offer much help in deducing causation from association, because there is little evidence to support the basic assumptions (Freedman [18]). To pursue the discussion here, a slightly more explicit definition of the models may be in order.

DEFINITION. A “path model” starts with variables at “level 0,” which are exogenous in the minimal sense that they are not explained within the model. Variables “at level 1” are built up as linear combinations of level 0 variables, plus independent random errors. More generally, variables “at level  $k$ ” are built up as linear combinations of variables at previous levels; again, there are additive, independent random errors. Variables at level 1, level 2,  $\dots$  are “endogenous,” in the sense that they are explained within the system. The path model may be presented as a “path diagram,” like Fig. 1, or Fig. 5 below. Nodes represent variables in the model; if there are arrows from  $X, Y, \dots$  to  $Z$ , then  $X, Y, \dots$  are explanatory variables in the regression equation for  $Z$ . Nodes are often called “vertices,” and the diagrams are referred to as “graphs” or “causal graphs.”<sup>15</sup>

The path model may represent mere association—conditional dependence and independence relations. Or the model may represent causation. I will take that up later. For now, however, either interpretation suffices.

<sup>14</sup> The model used by Rindfuss *et al.* [55] would not fall into this category, if ED and AGE really influenced each other. The SGS [62] framework excludes reciprocal causation, by assumption; so do path models, as I define them. However, some authors extend the definition of path models to include simultaneous equation models for reciprocal causation.

<sup>15</sup> SGS [62] seem to make the strong—and quite unusual—assumption that exogenous variables are independent of each other. That may be part of the reason why their algorithms estimate such peculiar models in Figs. 5 and 6 below. There is another, even more esoteric, point. To estimate an equation, its error term need only be assumed independent of the explanatory variables. If so, error terms from different equations may be correlated; then standard procedures for computing the correlations among the variables will not apply: see Freedman [18, pp. 112–114]; Seneta [57, p. 199]. SGS seem to interpret correlated errors as indicating the presence of “latent variables.” Such variables will be mentioned in notes to Figs. 5 and 6, below.



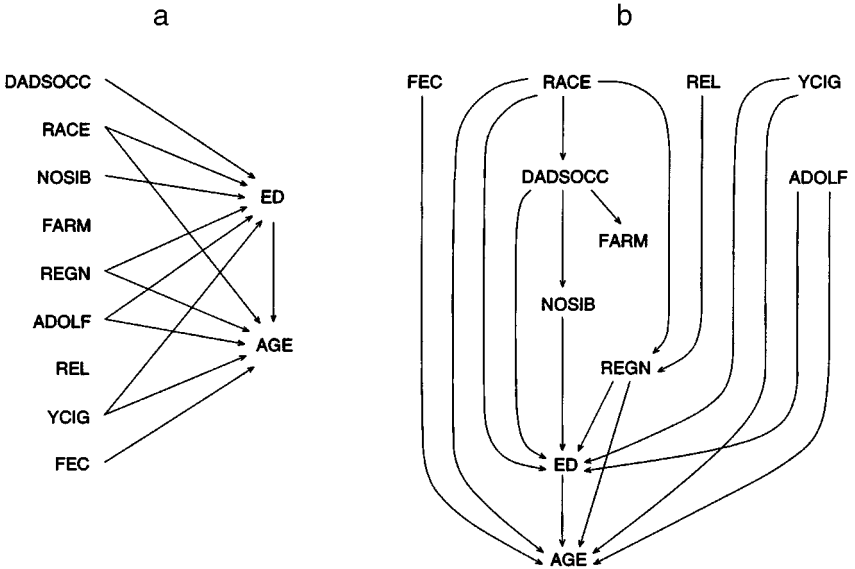


FIG. 5. The left-hand panel shows the model reported by SGS [62]. The right-hand panel also shows connections among the regressors, as determined by the search program TETRAD. BUILD indicates that latent variables are present, i.e., errors are correlated across equations. BUILD asks whether it should assume “causal sufficiency”; without this assumption [62, p. 45], the program output is uninformative. Therefore, I told BUILD to make the assumption; I believe that is what SGS [62] did for the Rindfuss example. Also see Spirtes *et al.* [63, pp. 13–15]. I told BUILD that ED and AGE could not cause the remaining variables, following SGS [62, p. 139]. However, SGS [62] actually made the stronger assumption that (i) FEC, ED, and AGE could not cause YCIG, and (ii) FEC, ED, AGE, and YCIG could not cause the remaining variables. With the assumption of causal sufficiency, BUILD seems to use the PC algorithm; without the assumption, the FCI algorithm comes into play. Much of this information comes from Richard Scheines (personal communication). Data are from Rindfuss *et al.* [55], not SGS [62]; with the SGS covariance matrix, FARM causes REGN and YCIG causes ADOLF.

Suppose the graph is “sparse”—each equation in the model involves relatively few variables. Suppose, too, there are no troublesome algebraic identities among the regression coefficients; in SGS terminology, the distribution is “faithful” to its graph [62, p. 35]; see Section 11.2 below. You have a sample—many independent realizations of variables  $X, Y, Z, \dots$ . You are willing to assume the distribution conforms to a path model, but do not know which model. You do not even know which variables are at level 0, which are at level 1, and so forth.

SGS [62] claim their algorithms are likely to find the underlying path model, or a rather similar model, and in short order. Their most convinc-

ing evidence is based on simulation experiments, where the computer generates data from a path model and the SGS algorithms try to infer the model from the data [62, pp. 145ff, 152ff, 250ff, 320ff, 332ff]; in these experiments, the algorithms do very well. Roughly speaking, the SGS algorithms are variants of “best subsets” regression, the search being over graphs rather than subsets. The data come into the SGS algorithms only through the covariance structure. The rest of the apparatus—the diagrams, the Markov property, faithfulness, etc.—consists of assumptions.

SGS [62] seem to assert that their algorithms determine causality, as a matter of mathematics. Such assertions are not defensible. In the SGS formalism, causation is obtained not by mathematical proof but by mathematical assumption. If you assume that the arrows in the underlying path diagram represent causes, then the arrows found by the algorithms represent causes. If you assume that the underlying arrows represent mere associations, then the arrows found by the algorithms represent associations. Causation has to do with empirical reality, not with mathematical proofs based on axioms. The issue is not one of theorems, but of the connection between theorems and reality.

The SGS algorithms [62], like many earlier statistical procedures (factor analysis, LISREL, etc.), proceed by analyzing the correlation matrix of a set of variables. I will call such methods “correlational.” Sections 7–10 consider applications of the SGS algorithms to real examples. Sections 11–12 try to explain the key ideas in the SGS formalism and indicate by mathematical example some of the intrinsic limitations. Before proceeding, however, I discuss the SGS statement of assumptions.

### 6.1. *The SGS Statement of Assumptions*

SGS [62] discuss the role of assumptions in their theory several times (pp. 53–69, pp. 75–81, pp. 324–325, p. 351). However, the clearest statement can be found when SGS are trying to discredit the evidence that smoking causes lung cancer:

effects \* \* \* \* cannot be predicted from \* \* \* \* sample conditional probabilities. [p. 302]

Readers may consult the original for context, to see whether the omitted material affects the meaning. The advantage of the quote is clarity. If the statement is generally applicable, then SGS—like Yule and Pearl before them—have disavowed the ability to infer causation from association.

## 7. THE SGS EXAMPLES

SGS [62] share my pessimistic views about regression. They claim, however, that their algorithms will succeed where regression has failed:

In the absence of very strong prior causal knowledge, multiple regression should not be used to select the variables that influence an outcome or criterion variable in data from uncontrolled studies. So far as we can tell, the popular automatic regression search procedures [like stepwise regression] should not be used at all in contexts where causal inferences are at stake. Such contexts require improved versions of algorithms like those described here to select those variables whose influence on an outcome can be reliably estimated by regression. In applications, the power of the specification searches against reasonable alternative explanations of the data is easy to determine by simulation . . . . [p. 257]

At first reading, SGS seems to be filled with real examples showing the successful application of their algorithms. That is an illusion. Many of the examples are based on simulation, and I set those aside.<sup>16</sup> The real examples are mostly to be found in SGS [62, pp. 132–152, 243–256].<sup>17</sup>

The main examples given in SGS [62] are path models. But these cannot withstand scrutiny—see Section 5 above and Sections 8–9 below. One exception is the stratification model of Blau and Duncan [3]. SGS [62, pp. 142–145] seem to be quite critical of this model; their current position is almost diametrically opposite to the one in Glymour *et al.* [24, pp. 33–39]. Like SGS, I do not believe that the Blau–Duncan regressions are a satisfactory causal model. On the other hand, as descriptions of the data, the equations can tell us something important about our society: see Freedman [18, pp. 122, 220]. The discussion in SGS adds little to our understanding either of the model or of stratification.

SGS [62] appear to use the health effects of smoking as a running example to illustrate their theory.<sup>18</sup> Again, there is an illusion. The causal diagrams are all hypotheticals, no contact is made with data, and no substantive conclusions are drawn. If the diagrams were proposed as real descriptions of causal mechanisms, they would not be credible.

What about the substantive question: does smoking cause lung cancer, heart disease, and many other illnesses? SGS [62] appear not to believe the

<sup>16</sup> Simulations tell us how well the SGS algorithms do *if* the underlying statistical assumptions hold good; the assumptions are built into the computer code that generates the simulated data. When applying statistical algorithms to real data, a critical question is *whether* those assumptions hold. The simulations do not address such questions.

<sup>17</sup> Parallel material is in [23, pp. 13–16, 21–23].

<sup>18</sup> See, e.g., [62, pp. 18, 216–237].

epidemiological evidence. When they actually get down to arguing their case, they use a rather old-fashioned method—a literature review with arguments in ordinary English [62, pp. 291–302]. Causal models and search algorithms have disappeared.

I approve of the method if not the implementation; the summary is wrong in some places and tendentious in others. However, the review does show the complexity of the issues. To make judgments about causation, you need to consider death certificate data, necropsy data, case control and cohort studies, twin studies, dose response curves, as well as animal experiments and human experiments. The force of the epidemiological evidence—and the SGS critique—depends on the complex interplay among these various studies and data sets.

In the end, SGS [62] do not really make bottom-line judgments on the health effects of smoking, at least so far as I can see. Their principal conclusion is methodological: nobody understood the issues.

Neither side understood what uncontrolled studies could and could not determine about causal relations and the effects of interventions. The statisticians pretended to an understanding of causality and correlation they did not have; the epidemiologists resorted to informal and often irrelevant criteria, appeals to plausibility, and in the worst case to *ad hominem* . . . . While the statisticians didn't get the connections between causality and probability right, the . . . . "epidemiological criteria for causality" were an intellectual disgrace, and the level of argument . . . was sometimes more worthy of literary critics than scientists. [62, pp. 301–302].

Part of a sentence in SGS [62, p. 4] does seem to grant one of the major claims made by the epidemiologists, "smoking does cause lung cancer." But that only complicates the puzzle. If you don't believe the evidence, why accept the claim?

Despite SGS [62], the epidemiologists did have a good understanding of the issues and made a strong case against smoking. The arguments were imperfect, and some reasonable doubts may remain. But the data, taken all in all, are compelling. The epidemiological literature on smoking is far stronger than anything I have seen in the social sciences. For a survey of the evidence, see Cornfield *et al.* [7]; this paper is still worth reading. More recent data are reviewed in [30].

SGS [62] elected not to use their analytical machinery on the smoking data—a remarkable omission. When applied to the examples that SGS actually chose, the algorithms produce one small disaster after another, as will now be seen. In sum, SGS [62] claim to have developed techniques for generating causal models; but they do not have any success stories.

## 8. USING THE SGS SEARCH PROCEDURE

The SGS search procedures are embodied in a computer program called TETRAD [62]. Version 2.1 of this program was kindly provided by Richard Scheines and Peter Spirtes. The BUILD module is the part of TETRAD used to discover path models with no latent variables. I ran BUILD on two examples—Rindfuss *et al.* [55] and AFQT (to be discussed in Section 9).

8.1. *Rindfuss et al.*

To explain AGE (age at first birth) in the Rindfuss *et al.* [55] example, the SGS [62] algorithms select the variables shown in Table III. Regression

TABLE III

The SGS [62] model for age at first birth, computed using the SGS covariance matrix or the Rindfuss *et al.* [55] covariance matrix

	SGS covariance			Rindfuss <i>et al.</i> covariance		
	$R^2 = 0.27$			$R^2 = 0.24$		
	Estimate	SE	<i>t</i>	Estimate	SE	<i>t</i>
RACE	-1.66	.30	-5.50	-1.66	.30	-5.46
REGN	-0.56	.19	-3.01	-0.63	.19	-3.35
ADOLF	1.89	.22	8.60	2.01	.22	8.98
YCIG	2.14	.25	8.63	-0.89	.25	-3.53
FEC	2.72	.28	9.70	2.77	.28	9.72
ED	0.67	.04	18.00	0.60	.04	15.72

*Note.* (i) Intercepts are not reported; OLS estimates.

(ii) The first column in Table 3 shows parameter estimates. The second shows standard errors, or SEs, which indicate the likely size of the differences between the estimates and the true parameter values. The *t*-statistics in the third column are the ratios of estimates to SEs. Generally, a *t*-statistic above 2 or 3 in absolute value indicates that the corresponding parameter is unlikely to be truly 0. For details, see the Appendix.

(iii) The parameters are features of the model, and the SEs are computed using the model. If you do not believe in the existence of the parameters apart from the data, or do not accept the statistical assumptions in the model, the SEs and *t*-statistics are likely to be meaningless. In any case, performing multiple tests—as in a search algorithm—complicates the interpretation of the *t*-statistics [17, 23].

(iv)  $R^2$  is generally interpretable as a descriptive statistic, whether or not the assumptions of the model hold true. An  $R^2$  of 0.27 indicates that about 27% of the variance in AGE has been explained; that is not much, and models in the social science literature often have even less explanatory power. For a discussion of  $R^2$ , see [20, pp. 78–81].

(v) According to current epidemiological opinion, smoking does have some biological effect, delaying conception by several weeks. However, the women who choose to smoke are different from the nonsmokers and have their first child almost a year earlier. This effect remains even after controlling for the measured background factors in the regression; the coefficient of YCIG is -0.89 years.

estimates for the coefficients, based on summary data in SGS, are reported in the first three columns of the table. The coefficients for ADOLF (the indicator for women from broken homes) and YCIG (an indicator for smoking by age 16) have positive signs. That is paradoxical: women from broken homes and women who smoke should be having children earlier, not later.<sup>19</sup> The signs should be negative, not positive. SGS do not comment on this issue.

Rindfuss *et al.* [55] give standard deviations and correlations for their data; SGS [62, p. 139] used these statistics to compute a covariance matrix, but reversed some of the signs. The last three columns of Table III report regression estimates computed from the correct covariances. The problem with YCIG disappears, but the sign for ADOLF stays positive. Anyone can make a mistake entering data; ignoring paradoxical signs in a causal model is quite another matter.

SGS [62] report only a graphical version of their model. They say,

Given the prior information that ED and AGE are not causes of the other variables, the PC algorithm (using the .05 significance level for tests) directly finds the model [in Figure 5(a)] where connections among the regressors are not pictured. [62, p. 139]

However, connections among regressors can be of interest. Although TETRAD is supposed to discover the causal ordering of explanatory variables, it produces the very strange model shown in Fig. 5(b). For example, the model says that race and religion cause region of residence. Comments on the sociology may be out of place, but consider the statistics. The equation is

$$\text{REGN} = a + b \times \text{RACE} + c \times \text{REL} + \epsilon. \quad (7)$$

REGN is a dummy variable, coded 1 for respondents who grew up in the South, 0 for others; RACE is 1 for black respondents and 0 for others; REL is 1 for Catholics, 0 for others;  $\epsilon$  is normally distributed. In consequence, this equation forces impossible values on REGN: the left-hand side is 0 or 1, the right-hand side varies from  $-\infty$  to  $+\infty$ . Now  $R^2$  is only 0.16, so  $\epsilon$  contributes most of the variance; Eq. (7) can hardly be defended as an approximation. Having dummy variables in the middle of path diagrams is a blunder. (FARM creates a similar problem; so does NOSIB, although less extreme.) In short, the SGS algorithms have produced a model that fails the most basic test—internal consistency.

<sup>19</sup> Smoking, broken homes, and early childbearing seem to be correlates of social disadvantage and indicators of personality traits. DADSOCC and RACE are quite imperfect controls for family background; therefore, YCIG and ADOLF are likely to pick up the effects of background, as well as the effects of omitted personality variables. See note (v) to Table III. This sort of bias is discussed in Section 12.2 below. Also see Clogg and Haritou [6].

## 9. THE ARMED FORCES QUALIFICATION TEST

SGS [62] discuss an example based on the Armed Forces Qualification Test (AFQT).<sup>20</sup> The AFQT is a linear combination with fixed weights of scores on certain subtests. Some of these subtests, as well as subtests that are not part of the AFQT, are listed in Table IV. The problem is to decide which subtests go into the AFQT and which do not.

The problem may be stated more algebraically as

$$\begin{aligned} \text{AFQT score} = & a_1 \times \text{NO} + a_2 \times \text{WK} + \cdots + a_7 \times \text{GS} \\ & + b_1 \times \text{UN}_1 + \cdots + b_n \times \text{UN}_n, \end{aligned} \quad (8)$$

where  $\text{UN}_1, \dots, \text{UN}_n$  are unobservable. Some of the  $a$ 's are zero, and the challenge is to figure out which ones.

We have data on 6224 subjects, summarized as a covariance matrix. According to SGS [62, pp. 243–244]:

a linear multiple regression of AFQT on the other seven variables gives significant regression coefficients to all seven and thus fails to distinguish the tests that are in fact linear components of AFQT. . . . Given the prior information that AFQT is not a cause of any of the other variables, the PC algorithm in TETRAD II correctly picks out {AR, NO, WK} as the only . . . variables that can be components of AFQT. . . .

To test the claims about regression, I ran AFQT on all the observable subtests. As Table V shows, EI and MC are related to AFQT only at the chance level. Moreover, MK and GS have negative coefficients, but psychometric practice frowns on subtests that are negatively related to overall test scores. It is a natural conjecture that NO, WK, and AR go into AFQT, while the other four subtests do not. Contrary to the claims of SGS, the AFQT can be handled by ordinary statistical methods.

TABLE IV  
Subtests Analyzed by SGS [62]

1. Numerical Operations	NO
2. Word Knowledge	WK
3. Arithmetical Reasoning	AR
4. Mathematical Knowledge	MK
5. Electronics Information	EI
6. Mechanical Comprehension	MC
7. General Science	GS

*Note.* Some go into the AFQT and some do not.

<sup>20</sup> SGS [62, p. 243]. Institutional background on the AFQT will be found in Section 12.5.

TABLE V  
Regression of AFQT on All the Observable Subtests

	Estimate	SE	<i>t</i>
NO	0.24	.022	10.8
WK	1.17	.029	40.5
AR	1.03	.028	36.4
MK	-0.24	.028	-8.7
EI	-0.03	.024	-1.3
MC	0.03	.024	1.3
GS	-0.13	.029	-4.6

*Note.* Variables were centered at their means.

The AFQT problem is in some ways quite easy. By definition, the “causes” or subtests combine linearly with the parameters to produce the AFQT as an “effect.” Joint normality of test scores seems to follow from the procedures used to construct the tests: consequently, scores on any one subtest can be presented as a linear combination of other subtest scores, with additive random errors. Thus, critical issues in most empirical studies have disappeared.<sup>21</sup>

### 9.1. TETRAD

According to SGS, given the prior information that AFQT does not cause the other variables, TETRAD correctly picks out AR, NO, and WK as the components of the AFQT.<sup>22</sup> Without that prior information, however, TETRAD declares AFQT to be the *cause* of these subtests, rather than the *effect*. With the prior information, TETRAD produces the strange results shown in Figure 6.<sup>23</sup> Now, for instance, the subtest NO may “cause” the overall test score AFQT, but it can hardly cause the other subtests AR or MK. Furthermore, there is a cycle in the figure:

$$MC \rightarrow AR \rightarrow WK \rightarrow GS \rightarrow MC.$$

In principle, such cycles were excluded by prior assumption, as well they might be. Subtests should not cause themselves, even indirectly. To sum up:

- (i) ordinary least squares techniques pick out NO, AR, and WK for the probable components of the AFQT, just as TETRAD does;
- (ii) TETRAD produces the curious model in Figure 6.

<sup>21</sup> On the other hand, unobserved variables may create serious problems (Section 12.4).

<sup>22</sup> SGS [62, p. 243].

<sup>23</sup> The program output is given in Spirtes *et al.* [63, pp. 10–11].



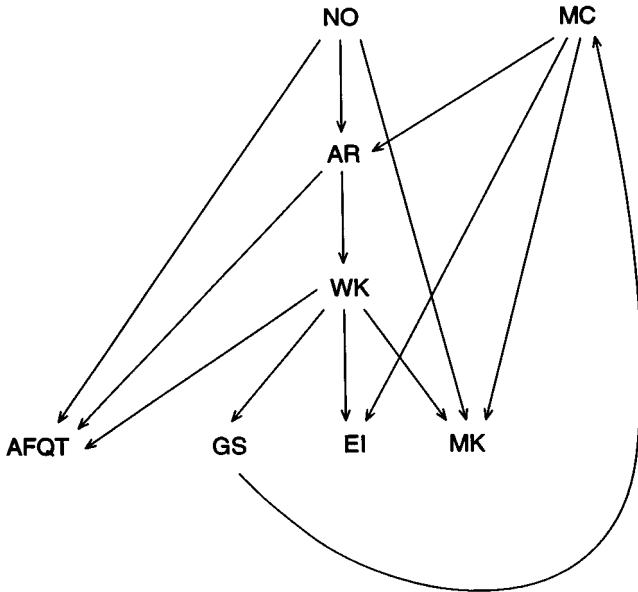


FIG. 6. AFQT and its subtests arranged in causal order by the search program TETRAD. I believe SGS [62, pp. 243–244] used BUILD, with the assumption of causal sufficiency, for the AFQT example. Also see Spirtes *et al.* [63, pp. 8–11]. The program indicates there are latent variables, i.e., correlations in the errors.

## 10. FOREIGN INVESTMENT AND POLITICAL OPPRESSION

As noted in Section 7, SGS are quite pessimistic about typical social-science applications of regression. While I agree with the bottom line, their specific objections seem misplaced. One example is enough to make the point. Timberlake and Williams [65] offer a regression model to explain political exclusion (PO) in terms of foreign investment (FI), energy development (EN), and civil liberties (CV). High values of PO correspond to authoritarian regimes that exclude most citizens from political participation; high values of CV indicate few civil liberties. Data come from 72 countries. Correlations among the Timberlake–Williams variables are shown in Table VI.

The equation proposed by Timberlake and Williams [65] is

$$PO = a + b \times FI + c \times EN + d \times CV + \text{error}. \quad (9)$$

Empirical results are shown in the first three columns of Table VII. The estimated coefficients of FI is significantly positive and is interpreted as

TABLE VI  
The Timberlake and Williams Correlation Matrix

	PO	FI	EN	CV
PO	1.000	-.175	-.480	.868
FI	-.175	1.000	.330	-.391
EN	-.480	.330	1.000	-.430
CV	.868	-.391	-.430	1.000

*Note.* Correlation matrix for political oppression (PO), foreign investment (FI), energy development (EN), and civil liberties (CV). Source: [62, p. 249].

measuring the effect of foreign investment on political exclusion: see Timberlake and Williams [65, p. 143].

SGS discuss this example [62, pp. 248–250], suggesting that Timberlake and Williams have confused cause and effect. The alternative causal sequence is not spelled out. Presumably, the idea is that dictators “cause” foreign investment in the sense that investors think dictatorial regimes offer greater stability, etc.

The main step in the SGS statistical argument comes down to this: the correlation of  $-0.175$  between political exclusion and foreign investment is at the chance level. The calculation rides on two assumptions: (i) the 72 countries in the data set are a random sample from some much larger set of countries and (ii) the variables follow a multivariate normal distribution. These time-honored but madcap assumptions are not stated explicitly by

TABLE VII  
The Timberlake and Williams Model

	$R^2 = .81$			$R^2 = .93$		
	Estimate	SE	$t$	Estimate	SE	$t$
FI	.23	.059	3.9	.44	.036	12
EN	-.18	.060	-2.9	-.22	.037	-6
CV	.88	.061	14.4	.95	.038	25

*Note.* Political exclusion (PO) is regressed on foreign investment (FI), energy development (EN), and civil liberties (CV). The first three columns show results for the observed correlation matrix (Table VI). The last three columns show what happens when  $r(\text{PO}, \text{FI})$  is set to 0. Coefficients in Table VII are standardized, that is, computed from variables standardized to have mean 0 and variance 1. The coefficients reported by SGS [62, p. 249] are not standardized and therefore do not match the correlation matrix.

SGS, let alone justified. (Of course, the assumptions behind the statistics in Timberlake and Williams might seem equally antic.)

However, for the sake of argument, let us grant SGS [62] their assumptions. On that basis, the standard error for the correlation in question is about  $1/\sqrt{72} \approx .12$ . I change the suspect correlation coefficient from its observed value of  $-0.175$  to the new value of  $0$ , a difference of about  $1.5$  SEs. I then recompute the model (last three columns in Table VII). The results are even better for Timberlake and Williams: the estimated coefficients are bigger and more significant; the signs stay the same; and  $R^2$  moves closer to  $1$ .<sup>24</sup>

I will not defend the model any further. Measurement problems are extreme, and the list of omitted variables very long. SGS may well be right, that cause and effect have been confused. But the demonstration is peculiar. The correlation matrix cannot show that FI, EN and CV cause PO—the fatal flaw in the Timberlake-Williams model. (Of course, Timberlake and Williams are not alone in this respect.) Nor can the matrix show that FI, EN and CV do not cause PO—the corresponding flaw in SGS. Indeed, it is trivial to construct four variables labelled FI, EN, CV and PO, such that FI, EN and CV do cause PO; but sample correlation matrices will look rather like the one in Table VI. This only sharpens the basic question. What do any of these calculations tell us about the world outside the computer?

## 11. SOME MATHEMATICAL ISSUES

Sections 11 and 12 address by mathematical example two questions:

- (i) To what extent can correlational methods recover an underlying path diagram?
- (ii) When can the arrows in the diagram be interpreted as indicating causation, rather than conditional independence and dependence?

The examples will indicate how SGS [62] use the “faithfulness” assumption to help them answer such questions. Issues of identifiability and consistency will be discussed, and methodological contributions in SGS will be delineated. Sections 11 and 12 are more technical than previous material; readers can skip to Section 13 without losing the thread of the argument.

The focus is on linear models. Suppose you have a covariance matrix that describes certain variables. Assume these variables are jointly normal,

<sup>24</sup> The new matrix is still positive definite, so it is a legitimate correlation matrix. Section 12.1 discusses the connection between the Timberlake-Williams model and the faithfulness assumption. Also see Cartwright [5, pp. 79–84].

with mean 0; that avoids all questions of linearity etc. and all problems created by having only finite amounts of data. However, the statistical procedures I am considering—like the SGS algorithms—will operate on that covariance matrix and on nothing else. Such procedures may be called “correlational.”

Path models were defined in Section 6. Briefly, you start with variables at level 0; variables at level  $k$  are linear combinations of variables at lower levels, plus independent random errors. In a path diagram, nodes represent variables. There is an arrow from  $X$  to  $Y$  if  $X$  is used as an explanatory variable in the equation for  $Y$ .

Exogeneity is a critical concept. As indicated before, the term is used in at least three senses. The weakest definition is purely mechanical: exogenous variables are not explained within the model, but are supplied to the model. Variables at level 0 in a path model are exogenous in this minimal sense. A more restrictive definition: exogenous variables are statistically independent of the error terms in the equations. The third idea is the one that is relevant to causal inference:  $X$  is exogenous if selecting subjects with  $X = x$  gives the same results as intervening to set  $X = x$ .

There are tests for exogeneity in the literature, as well as model specification tests. However, these have limited relevance to causal inference. For example, Hausman [26] assumes that certain variables are known *a priori* to be exogenous and then tests whether other variables are exogenous; he interprets exogeneity as orthogonality to disturbance terms. He also has a test that detects correlation between errors from equations in a path model. White [69, 70] focuses on similar issues—for instance, testing whether the variables have a jointly normal distribution.

Another reference in the econometric literature is Engle, Hendry, and Richard [15]. These authors distinguish several kinds of exogeneity; “strict” exogeneity means independence of variables and error terms, but only “super” exogeneity permits estimating the effects of interventions. Examples are given to illustrate the definitions [15, pp. 287–294]. There is further discussion in Leamer [35].

### 11.1. *The Basic Statistical Problem*

Suppose you have  $n$  random variables with a jointly normal distribution; all the variables have mean 0, and you know the covariance matrix, which is positive definite. You wish to present this covariance matrix as a path model. In a sense, nothing is easier. Simply order the variables, arbitrarily, as  $X_1, X_2, \dots, X_n$ . By successively applying regression, we can find coefficients  $a_{ij}$  and error terms  $\epsilon_i$ , such that  $X_1, \epsilon_2, \dots, \epsilon_n$  are all independent

with mean 0, and Eq. (10) holds:

$$\begin{aligned} X_2 &= a_{21}X_1 + \epsilon_2 \\ X_3 &= a_{31}X_1 + a_{32}X_2 + \epsilon_3 \\ &\vdots \\ X_n &= a_{n1}X_1 + \cdots + a_{n,n-1}X_{n-1} + \epsilon_n. \end{aligned} \quad (10)$$

Then  $X_1$  is presented as exogenous and the “cause” of  $X_2$ ; next,  $X_1$  and  $X_2$  “cause”  $X_3$ ; and so forth. In short, there are many ways to present a covariance matrix as a path diagram; few if any will be relevant for causal inference.<sup>25</sup>

### 11.2. The Faithfulness Assumption

How can you single out one path diagram from the many that correspond to a given covariance matrix? At this point, SGS [62] seem to use the “faithfulness” assumption; this assumption is also used to handle confounding, as discussed in Section 12.1 below. Basically, a covariance matrix is faithful to a diagram provided conditional dependencies and independencies are determined by the presence or absence of arrows in the diagram, rather than specific numerical values of parameters.

By way of example, Fig. 7 shows two path diagrams. On the left,  $X$  causes  $W$  through the intervening variables  $Y$  and  $Z$ ; on the right, the flow of causality is reversed.<sup>26</sup> The lower case letters on the arrows stand for “path coefficients,” that is, standardized regression coefficients. How could SGS [62] distinguish between the two theories in the figure? Their idea seems to be as follows:

In the left hand diagram,  $Y$  and  $Z$  are conditionally independent given  $X$ ; on the right, however,  $Y$  and  $Z$  are conditionally dependent given  $X$ .

Another contrast:

In the left-hand diagram,  $Y$  and  $Z$  are conditionally dependent given  $W$ ; on the right, however,  $Y$  and  $Z$  are conditionally independent given  $W$ .

<sup>25</sup> For the construction in (10), simply choose  $a_{21}$  so  $E\{X_2|X_1\} = a_{21}X_1$ ; choose  $a_{31}$  and  $a_{32}$  so  $E\{X_3|X_1, X_2\} = a_{31}X_1 + a_{32}X_2$ ; and so forth. For details, see the Appendix. Since the ordering of the variables in (10) is arbitrary, fitting such equations or drawing path diagrams cannot determine which variables are causes and which are effects. In particular,  $X_1$  may be exogenous in the sense that it is statistically independent of disturbance terms; that by itself does not suffice to estimate the results of manipulating  $X_1$ , since we cannot tell whether  $X_1$  is a cause or an effect.

<sup>26</sup> In this section, I use “cause” in its ordinary (perhaps undefinable) sense. However, the technical point—about the possibility of estimating path diagrams from covariance matrices—still holds if the arrows are interpreted as merely representing association. “Causation” is then colorful shorthand (perhaps too colorful) for a certain kind of covariation.

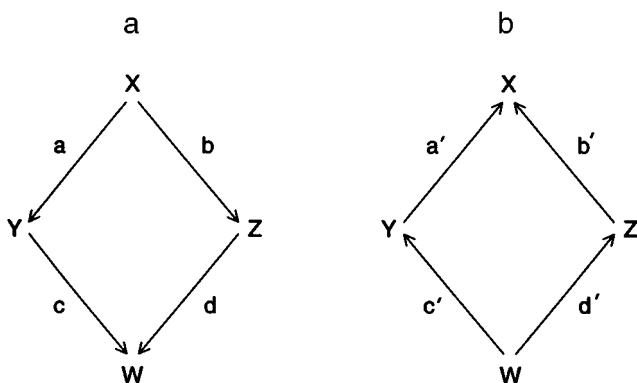


FIG. 7. If two path diagrams have the same covariance matrix, correlational methods cannot tell them apart; the faithfulness assumption is made to rule out such problems. The lower case letters on the arrows denote “path coefficients,” that is, standardized regression coefficients.

Therefore, the pattern of conditional dependence and independence identifies the diagram. (In both diagrams,  $X$  and  $W$  are conditionally independent given  $Y$  and  $Z$ .)

This idea works for many path diagrams, but fails for others. Indeed, the path coefficients can be chosen so the pattern of conditional dependence and independence is the same in the two diagrams. Even worse, both diagrams can give rise to the same covariance matrix—so correlational methods cannot tell which is right. SGS [62] make the “faithfulness assumption” in order to rule out such indeterminacies. (The workings of the assumption will be explained below.)

However, that only moves the difficulty to another place. Faithfulness is hardly an empirical fact; it is an assumption about unobservables, made to rule out situations that cannot be handled by correctional methods. The SGS analytical program can now be stated rather simply. If the arrows in a path diagram represent causation not association, and if the path diagram can be estimated from data, then SGS can indeed infer causation from association.

The balance of Section 11.2 provides technical backup; readers can skip to Section 11.3. The left-hand panel in Fig. 7 is described by

$$Y = aX + \delta_1, \quad Z = bX + \delta_2, \quad W = cY + dZ + \delta_3. \quad (11)$$

In this equation,  $X, \delta_1, \delta_2, \delta_3$  are independent and normal, with mean 0;  $X, Y, Z, W$  all have variance 1. The covariance matrix of  $X, Y, Z, W$  can be

computed from the four parameters  $a, b, c, d$  as shown in (12):

	$X$	$Y$	$Z$	$W$
$X$	1	$a$	$b$	$ac + bd$
$Y$	$a$	1	$ab$	$c + abd$
$Z$	$b$	$ab$	1	$d + abc$
$W$	$ac + bd$	$c + abd$	$d + abc$	1

(12)

It is a little theorem, which follows by a tedious calculation from (48) in the Appendix, that

$$\text{cov}(X, W|Y, Z) = 0. \quad (13)$$

This is an example of a conditional independence relation forced by a graph; (13) holds whatever the path coefficients in Fig. 7 may be.

The diagram on the left in Fig. 7 is reversible, provided

$$\text{cov}(Y, Z|W) = 0. \quad (14)$$

By (48) below, Eq. (14) is equivalent to

$$\text{cov}(Y, Z) = \text{cov}(Y, W) \times \text{cov}(Z, W). \quad (15)$$

By (12), this means

$$ab = (c + abd)(d + abc). \quad (16)$$

Rearranging (16) gives the quadratic equation

$$cd(ab)^2 - (1 - c^2 - d^2)ab + cd = 0. \quad (17)$$

One solution to (17) is

$$ab = \frac{1 - c^2 - d^2 - \sqrt{(1 - c^2 - d^2)^2 - 4c^2d^2}}{2cd}. \quad (18)$$

I chose  $a, c, d$  more or less at random, getting 0.1925, 0.2873, and 0.1245, respectively.<sup>27</sup> I computed  $b$  from (18), getting 0.2063. This choice forces the conditional independence relation (14) and violates the faithfulness assumption; conditional independence comes from the parameter values, not the presence or absence of arrows.

<sup>27</sup> There was a bit of luck here, because some values for  $a, c, d$  will not produce correlation matrices.

Given the values for the four parameters  $a, b, c, d$ , the covariance matrix (12) can be evaluated as

$$\begin{pmatrix} 1.0000 & 0.1925 & 0.2063 & 0.0810 \\ 0.1925 & 1.0000 & 0.0397 & 0.2922 \\ 0.2063 & 0.0397 & 1.0000 & 0.1359 \\ 0.0810 & 0.2922 & 0.1359 & 1.0000 \end{pmatrix}. \quad (19)$$

The path coefficients in the right-hand panel of Fig. 7 are easily computed from (19):

the path coefficient from  $W$  to  $Y$  is  $c' = \text{cov}(Y, W) = 0.2922$ ;

the path coefficient from  $W$  to  $Z$  is  $d' = \text{cov}(Z, W) = 0.1359$ ;

the path coefficients from  $Y$  and  $Z$  to  $X$  are obtained by multiple regression, as  $a' = 0.1846$  and  $b' = 0.1990$ .

With these choices, faithfulness does not hold and (19) can be represented by either diagram in Fig. 7. (For details on multiple regression, see the Appendix.) In effect, the faithfulness assumption precludes certain algebraic identities among the parameters, like (16). Since parameters are not observable, the faithfulness assumption is not subject to direct empirical tests based on finite amounts of data.

### 11.3. Complete Graphs

Even if the covariance matrix is faithful to a graph, however, problems of indeterminacy remain—particularly if the graph is “complete” in the sense that every pair of vertices is joined by an arrow. Figure 8 illustrates this indeterminacy. The same covariance matrix (20) for the variables

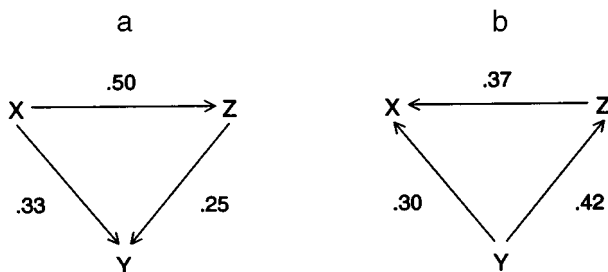


FIG. 8. Graphs (a) and (b) have the same covariance matrix. Both are complete; there is an arrow from every variable to every other variable. The numbers on the arrows are path coefficients, that is, standardized regression coefficients.



$X, Y, Z$  is represented either by the diagram in panel (a) or the one in panel (b), where the flow of “causality” is reversed:

	$X$	$Y$	$Z$
$X$	1	.46	.50
$Y$	.46	1	.42
$Z$	.50	.42	1

(20)

For a second example of indeterminacy when the graph is complete, consider four variables  $X, Y, Z, W$  with covariance matrix  $\Sigma$  given by

$$\Sigma = \begin{pmatrix} 1 & \frac{3}{4} & \frac{3}{4} & \frac{3}{4} \\ \frac{3}{4} & 1 & \frac{3}{4} & \frac{3}{4} \\ \frac{3}{4} & \frac{3}{4} & 1 & \frac{3}{4} \\ \frac{3}{4} & \frac{3}{4} & \frac{3}{4} & 1 \end{pmatrix}. \quad (21)$$

Figure 9 shows two complete path diagrams, both of which are compatible with the given covariance matrix. In the left-hand panel,  $X$  is exogenous and “causes”  $Y$ ; then  $X$  and  $Y$  “cause”  $Z$ ; finally,  $X, Y, Z$  “cause”  $W$ . In panel (b), the flow of “causality” is reversed. The equations corresponding to the left-hand panel are given as (22); panel (b) is described in (23):

$$\begin{aligned} Y &= \frac{3}{4}X + \delta_1 \\ Z &= \frac{3}{7}X + \frac{3}{7}Y + \delta_2 \\ W &= \frac{3}{10}X + \frac{3}{10}Y + \frac{3}{10}Z + \delta_3; \end{aligned} \quad (22)$$

$$\begin{aligned} Z &= \frac{3}{4}W + \epsilon_1 \\ Y &= \frac{3}{7}Z + \frac{3}{7}W + \epsilon_2 \\ X &= \frac{3}{10}Y + \frac{3}{10}Z + \frac{3}{10}W + \epsilon_3. \end{aligned} \quad (23)$$

The covariance matrix  $\Sigma$  is also compatible with the factor analysis model (24), where the unobservable exogenous variable  $U$  causes all four observables (right-hand panel of Fig. 9):

$$X = U + \zeta_1, \quad Y = U + \zeta_2, \quad Z = U + \zeta_3, \quad W = U + \zeta_4. \quad (24)$$

In each system of Eqs. (22)–(24), the error terms are assumed to be independent and normally distributed with mean 0; error terms are independent of the exogenous variable. As a technical matter, the covariance matrix (20) is faithfully represented by both graphs in Fig. 8. Likewise, the

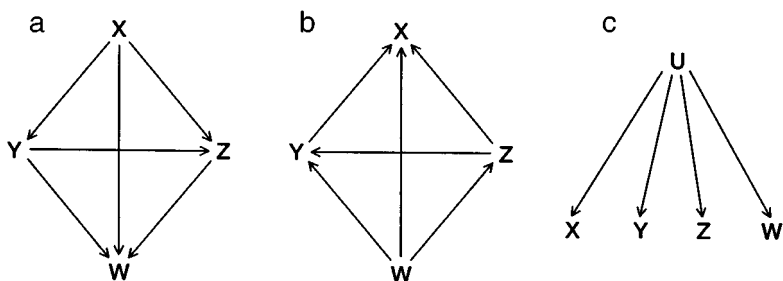


FIG. 9. Two complete path diagrams and a factor analysis model, all having the same covariance matrix.

covariance matrix (21) is faithful to Fig. 9(a) and to 9(b). Proofs may be based on (48) below.

To sum up, if a covariance matrix is faithful to a complete graph (with all pairs of vertices joined by arrows), it is faithful to many such graphs. Then correlational methods cannot tell the causes from the effects. SGS [62] techniques work best when the graph is sparse; that is, relatively few pairs of vertices are joined by arrows (Section 6).

#### 11.4. Identifiability and Consistency

The focus continues to be on linear models. In statistical terminology, models are “identifiable” when they make different predictions about observables. For example, suppose you have two models for your data. If, for all data sets,

$$P(\text{data}|\text{model 1}) = P(\text{data}|\text{model 2}),$$

there is an obvious problem—the data cannot distinguish between the models. If a path model is complete, or the faithfulness assumption is not imposed, then the graph underlying a covariance matrix is not identifiable; that is, the message of Sections 11.1–11.3. By way of illustration, the models in Fig. 7 are identifiable only if faithfulness holds.

However, even if we assume that a covariance matrix is faithful to a graph that is not complete, there may be several such graphs [62, p. 89]. For example, the following three graphs can generate the same covariance matrix:

$$X \rightarrow Y \rightarrow Z, \quad X \leftarrow Y \rightarrow Z, \quad X \leftarrow Y \leftarrow Z.$$

Thus, SGS do not seem to have succeeded in defining a class of graphs and covariance matrices for which identifiability holds [62, p. 194].

In statistical terminology, estimators are “consistent,” provided that, as the sample gets larger and larger, these estimators come closer and closer to the population parameters. If the parameters are not identifiable, however, consistency is problematic.

SGS [62] seem to claim that their algorithms will find all the path diagrams compatible with a given covariance matrix. However, the theorems suggest that the algorithms will at best find one such graph. SGS also seem to claim that their algorithms are consistent. However, without an identifiability theory for linear models, they cannot really be talking about consistency.

Statisticians do have the weaker notion of “Fisher consistency,” named after R. A. Fisher: when applied to data for the whole population, an estimator should reproduce the population parameters exactly. Theorems like 5.1 in SGS [62, p. 405] seem to demonstrate the analog of Fisher consistency, rather than anything stronger. Such theorems show that, given the population covariance matrix, the algorithms will produce one graph consistent with that matrix.

### 11.5. *Methodological Contributions*

There is a connection between the theory of “directed acyclic graphs” (DAGs) and the conditional independence of random variables. (See Darroch *et al.* [10], Kiiveri and Speed [34, 61], Pearl [43, 44], Verma and Pearl [49, 67], Geiger [22].) Much of this work is reviewed in SGS [62]. However, the mathematics of nonlinear causal diagrams seems to be irrelevant to the big question: how do we infer causation from association?

Most the applications in SGS are linear, i.e., based on path models. The “nonlinear causal diagrams” turn out to be multinomial models for categorical data; examples are in [62, pp. 147–151]. The issues about causation are quite similar to those for linear models, although the technical details are different.

This section will focus on path models. To describe the novelty in the SGS approach to estimation, suppose you have data from a path model and wish to estimate the model. Consider two cases:

*Case I.* You know the classification of variables as to level; that is, you know which variables are at level 0, which are at level 1, and so forth.

*Case II.* You do not know the classification of variables as to level.

In Case I, SGS [62] have little to tell us about estimation; as to confounding, see Section 12.1. Some of their algorithms seem to be equivalent to regression; others may be less efficient. In Case II, SGS try

to estimate the classification of variables as well as the path coefficients. That is the methodological contribution. To estimate the classification, SGS must impose the faithfulness assumption (Section 11.2). It is disappointing that SGS do not pin down the sense in which their algorithms are successful (Section 11.4).

## 12. MORE EXAMPLES AND SOME THEORY

Section 12.1 explains how the faithfulness assumption and conditional independence are supposed to eliminate confounding. Section 12.2 discusses omitted variables. Sections 12.3–12.5 revisit two examples from a more mathematical perspective; the idea is to show the limits of correlational methods.

### 12.1. *Faithfulness, Conditional Independence, and Confounding*

The problems created by unobservable variables are well known. As indicated above, SGS [62] handle such problems by imposing the faithfulness assumption. More specifically, the assumption is used to rule out confounding. If confounding can be eliminated, the goal is in sight—association may soon be converted into causation. This section, which is based on work by Jamie Robins (personal communication), examines the logic in more detail. Also see Pearl and Verma [49].

With some models, exact conditional independence forces a choice:

- either there is no confounding by unmeasured common causes,
- or the faithfulness assumption is violated.

Near-independence is not good enough; associations may then be entirely spurious. Thus, causal inferences made by the SGS technique need exact conditional independence as well as the faithfulness assumption.

This use of the faithfulness assumption has some theoretical interest. However, in order to base empirical work on such mathematical ideas, it would seem necessary to resolve the following questions, which SGS have not addressed:

- Can the basic models be validated?
- Can exact conditional independence be demonstrated?
- Given exact independence, why is exact cancellation of confounded effects overwhelmingly less likely than the total absence of such effects?

As a practical matter, exact independence seems quite unusual. However, the theory is worth understanding, and an example will make the

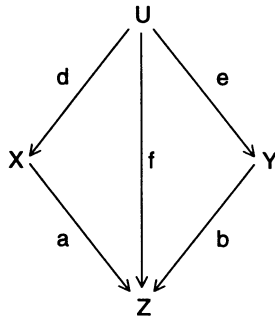


FIG. 10. The faithfulness assumption, conditional independence, and confounding. Variables  $X, Y, Z$  are observable;  $U$  is unobservable. Arrows represent causation, not just association. The lower-case letters on the arrows denote path coefficients. If a path coefficient vanishes, the corresponding arrow must be deleted.

position clearer. Figure 10 shows a relatively simple diagram where faithfulness and conditional independence would eliminate confounding. The arrows denote causation, not mere association. Variables  $X, Y, Z$  are observable;  $U$  is unobservable. Such unobservables are also called “confounders” or “unmeasured common causes.” The joint distribution is normal, and variables are standardized to have mean 0 and variance 1. The covariance matrix for all four variables is shown in (25).<sup>28</sup>

	$U$	$X$	$Y$	$Z$
$U$	1			
$X$	$d$	1		
$Y$	$e$	$de$	1	
$Z$	$f + ad + be$	$a + bde + fd$	$b + ade + fe$	1

(25)

Of course, only the covariance matrix (26) of the observables ( $X, Y, Z$ ) can be estimated from the data. In particular,  $de$  is determined from the observables, as  $\text{cov}(X, Y)$ :

	$X$	$Y$	$Z$
$X$	1		
$Y$	$de$	1	
$Z$	$a + bde + fd$	$b + ade + fe$	1

(26)

It may help to review the idea of faithfulness, in the context of our example. Faithfulness is an assumption about unobservables; more specifi-

<sup>28</sup> Covariance matrices are symmetric; only the lower triangular part is shown. Entries are assumed to be positive but less than 1. The matrix is assumed to be positive definite.

cally, it is a constraint on the relationship between the full covariance matrix (25) and the graph in Fig. 10. The assumption amounts to this: independence relationships (conditional and unconditional) are determined by the presence or absence of arrows in the diagram, not specific parameter values.

In particular, if the covariance matrix (25) is faithful to the diagram in Fig. 10, you cannot set any of the path coefficients to 0, except by deleting the corresponding arrow. An arrow from  $X$  to  $Z$ , say, entails that  $X$  has some causal effect on  $Z$ , no matter how small that effect may turn out to be.

I return to more conventional issues. In our example, the parameter of interest is  $b$ , the causal effect of  $Y$  on  $Z$ . Due to the unmeasured confounder  $U$ , a regression of  $Z$  on  $X$  and  $Y$  produces a biased estimate of  $b$ . By a slightly tedious calculation, the coefficient of  $Y$  in the regression equation is

$$b + fe(1 - d^2)/(1 - d^2e^2). \quad (27)$$

(For details on multiple regression, see the Appendix.) The bias in the regression estimate is the second term in (27). From a slightly different perspective,  $\text{cov}(Y, Z)$  in (26) measures the total association between  $Y$  and  $Z$ . Part of this association is real:  $b$  measures the causal effect of  $Y$  on  $Z$ . Alas, part of the association is spurious:  $ade + fe$  represents the effects of the confounder  $U$ .

The goal is to separate the real part of the association from the spurious part. The familiar obstacle is that we have only (26), not (25). And (26) does not suffice to separate  $b + ade + fe$  into its components. But, SGS might say, suppose that  $X$  and  $Z$  are conditionally independent given  $Y$ :

$$\text{cov}(X, Z|Y) = 0. \quad (28)$$

By (48) below, this means

$$\text{cov}(X, Z) = \text{cov}(X, Y) \times \text{cov}(Y, Z). \quad (29)$$

A bit of algebra based on (25) shows that (29) is equivalent to

$$a(1 - d^2e^2) + df = de^2f. \quad (30)$$

Although  $de$  is known and  $0 < de < 1$ , there are many possible ways to solve Eq. (30). At this point, SGS would invoke the faithfulness assumption, concluding that

$$a = 0, \quad f = 0. \quad (31)$$

The implication is that we have to remove the arrow from  $X$  to  $Z$ , as well as the arrow from  $U$  to  $Z$ .

Confounding has now been eliminated. On this basis,  $\text{cov}(Y, Z) = b$ , the whole of the association is real, and regression produces an unbiased estimate for the causal effect of  $Y$  on  $Z$ . At last, association has been converted into causation. (Of course, quite a lot of causality was built into Fig. 10 from the beginning—by assumption.)

Those were the implications of exact conditional independence. On the other hand, suppose we have approximate conditional independence:  $\text{cov}(X, Z|Y) = .00001$ . Now the faithfulness assumption has no force. Given the covariances in (26), we can match them by suitable choice of the other parameters, even if  $a = b = 0$ .<sup>29</sup>

With approximate conditional independence, observed associations can be entirely spurious. Thus, even in the realm of mathematics, faithfulness and conditional independence preclude confounding only when the independence is exact. To make the contrast sharper, let us assume faithfulness:

- If  $\text{cov}(X, Z|Y) = 0$ , then the association between  $Y$  and  $Z$  is purely causal; the effects of the unmeasured common cause  $U$  do not confound the relationship between  $Y$  and  $Z$ .

- If  $\text{cov}(X, Z|Y) = .00001$ , then confounding by unmeasured common causes may account for all of the observed association between  $Y$  and  $Z$ .

A similar problem must be considered when estimating path models from data (Section 11). Exact conditional independence, together with the faithfulness assumption, often permits us to identify the path diagram from the covariance matrix. However, approximate conditional independence is not enough; then, the covariance matrix will be faithful to a variety of complete graphs.

A final example is the Timberlake–Williams model (Section 10). This model explains political exclusion (PO) in terms of foreign investment (FI), energy development (EN), and civil liberties (CV); the sample correlation matrix was shown in Table VI. Consider three scenarios for the “true” correlation matrix  $\rho$ :

(i) Suppose  $\rho$  happens to equal the sample correlation matrix. Then, faithfulness obtains.

<sup>29</sup> This matching assumes, for instance, that any two of the variables have positive covariance, given the third. To avoid violating the faithfulness assumption, if you set  $a$  and  $b$  to 0, erase the corresponding arrows; if that is distasteful, set  $a$  and  $b$  to small but positive values. The SGS logic would apply to a wide variety of diagrams; however, an arrow from  $Y$  to  $X$ , no matter how small the coefficient, spoils the show.

(ii) Suppose the true correlation  $\rho(\text{PO}, \text{FI})$  between foreign investment and political exclusion happens to vanish exactly. Then, the Timberlake–Williams model violates the faithfulness condition; presumably, that is SGS’s real complaint.

(iii) If  $\rho(\text{PO}, \text{FI}) = .00001$ , faithfulness is restored. According to the SGS criteria, Timberlake and Williams are back in business.

Within the framework of path models, scenario (ii) cannot be rejected at conventional significance levels; neither can (iii); and (i) represents our best estimate, subject to large uncertainties. SGS seize on hypothesis (ii), the only one that legitimates their critique. They are balking at shadows.

## 12.2. *Omitted Variables*

The problem of omitted variables was raised by Cliff Clogg at the Notre Dame conference, and this section paraphrases one of his points. There is a response variable  $Y$ , with explanatory variables  $X$  and  $Z$ ; these may be construed as vectors. Suppose the data are generated according to the “true” model (32T):

$$Y = X\beta_R + \delta. \quad (32T)$$

The parameter vectors  $\beta$  and  $\gamma$  are unknown and to be estimated from data by regression; it is  $\beta$  that is of primary interest. Subjects are assumed to be independent and identically distributed;  $(X, Z)$  and the error term  $\epsilon$  are independent and jointly normal; all variables have expected value 0. Consider, too, the “restricted” model (32R), where  $\beta_R$  is defined so that  $E\{Y|X\} = X\beta_R$ ; the constituents of (32R) may be computed from the true model.<sup>30</sup>

$$Y = X\beta + Z\gamma + \epsilon. \quad (32R)$$

In principle, the variables  $X$ ,  $Y$ , and  $Z$  are all observable;  $X$  and  $Z$  may be correlated. However, investigators who do not know that  $Z$  is relevant may fit the restricted model R rather than the true model T. If so, the estimate of  $\beta$  can be quite biased. In the vernacular,  $\beta_R$  includes the effect of  $X$  on  $Y$  through  $Z$ . The covariance matrix of  $(X, Y)$  cannot distinguish between the two models, because the matrix can be generated

<sup>30</sup> Indeed,  $\beta_R = \beta + \alpha$ , where  $\alpha$  is obtained by the regression of  $Z\gamma$  on  $X$ . In other terms,  $Z\gamma = X\alpha + \eta$ , where  $\eta$  is normal with mean 0, independent of  $X$ . Then  $\delta = \epsilon + \eta$ . It may be seen that  $\alpha$  depends linearly on  $\gamma$ .



by either model. Therefore, no statistical procedure based on that matrix can tell you whether the restricted model is right or wrong.<sup>31</sup>

### 12.3. *On the Direction of Causality*

This section uses “cause” in its ordinary (perhaps undefinable) meaning, not as shorthand for certain kinds of covariation. I return to Judea Pearl’s example, shown in Fig. 2(a). Given the covariance matrix for  $X$ ,  $Y$ , and  $Z$ , the SGS [62] algorithm will produce the graph shown in panel (a). If you tell the algorithm that omitted variables are a possibility, it will tell you that  $Y$  cannot cause  $X$  or  $Z$ .

In the example,  $X$ ,  $Y$ , and  $Z$  are the only observables, and their covariance matrix is faithful to the graph in Fig. 2(a). I claim that such information cannot by itself determine the direction of the causal flow. To substantiate this claim, I now construct two theories. In both, the observables  $X$ ,  $Y$ , and  $Z$  will have the same covariance matrix, faithful to the graph in Fig. 2(a). However, the direction of the causal flow will be different in the two theories.

**THEORY 1.** I first generate  $X, Z, U$  as independent  $N(0, 1)$  variables;  $U$  is an unobservable error term. (If you want to intervene and change  $X$  or  $Z$ , now is your moment.) Then

$$Y = X + Z + U. \quad (33)$$

According to Theory 1,  $X$  and  $Z$  cause  $Y$ , as suggested by Fig. 2(a).

**THEORY 2.** I first generate  $Y$  as  $N(0, 3)$ . (If you want to intervene and change  $Y$ , now is your moment.) After a suitable pause, so that time’s arrow will delineate the flow of causality, I generate the errors  $V_1, V_2$ , and  $V_3$  as independent  $N(0, \frac{1}{3})$  variables and then produce  $X, Z$ , and  $U$ , according to

$$\begin{aligned} X &= \frac{1}{3}Y + V_1 - V_2 \\ Z &= \frac{1}{3}Y + V_2 - V_3 \\ U &= \frac{1}{3}Y + V_3 - V_1. \end{aligned} \quad (34)$$

In the second theory,  $Y$  causes  $X$  and  $Z$ . As far as the observables are concerned—namely, the joint distribution of  $X, Y$ , and  $Z$ —Theories 1 and 2 agree. Furthermore, the joint distribution is faithful to the graph in

<sup>31</sup> See Clogg and Haritou [6], who make the following very interesting point. Adding variables that are correlated with  $\epsilon$  can also bias the estimate of  $\beta$ ; this “included variable” bias can be just as troublesome as the more familiar “omitted variable” bias; the latter problem cannot be solved by throwing variables into the model. The SGS treatment of omitted variables was discussed in Section 12.1.

Fig. 2(a). But the direction of causality is determined neither by the data nor by the mathematics. With correlational methods, causality follows from the assumptions about the unobservables.

### 12.3. *The AFQT Problem*

SGS [62, p. 242] seem to claim that, as a demonstrable mathematical fact, their procedures will find the right answers:

Assuming the right variables have been measured, there is a straightforward solution to these problems: apply the PC, FCI, or other reliable algorithm, and appropriate theorems from the preceding chapters, to determine which  $X$  variables influence the outcome  $Y$ , which do not, and for which the question cannot be answered . . . . Then estimate the dependencies by whatever methods seem appropriate and apply the results of the previous chapter to obtain predictions of the effect of manipulating the  $X$  variables. No extra theory is required. We will give a number of illustrations . . . .

The first example given by SGS to illustrate this claim is AFQT (Section 9 above). To demonstrate that SGS are exaggerating more than a little, I pose a sharp mathematical question with the essential features of the AFQT problem. Then, I show the question to be undecidable by correlational methods. (Of course, when applied to the real example, both SGS and ordinary least squares made the right guess.)

To set up the question, assume that  $X$  and  $Y$  are random variables:  $X$  is a vector;  $Y$  is scalar.

$Y$  is a linear combination of  $X$ 's, with fixed weights. (35)

The observables are  $Y$  and  $V_1, \dots, V_7$ . (36)

Some  $V$ 's are  $X$ 's; some  $V$ 's are ringers. (A "ringer" is a variable that does not enter into the linear combination for  $Y$ .) There are also unobservables, including the  $X$ 's that are not  $V$ 's. Assume too that

The full joint distribution is multivariate normal, with mean 0. (37)

You are given the covariance matrix for the observables, but not the full covariance matrix. The problem is to say which of the  $V$ 's are  $X$ 's and which are ringers. I claim this problem is not solvable, because I can produce two different theories leading to different classifications of the  $V$ 's, but having the same joint distribution for the observables.

**THEORY 1.** I use the covariance matrix for the seven observable subtests  $V_1 = \text{NO}, \dots, V_7 = \text{GS}$ , together with the three unobservable subtests, CS, AS, and PC. (The subtests are listed in Table VIII, Section 12.5 below.) The full distribution is defined to be jointly normal, and all

variables have mean 0. Let  $Y = .5 \times \text{NO} + \text{AR} + \text{WK} + \text{PC}$ , where NO, AR, and WK are observable but PC is unobservable. In this theory,  $V_1, V_2, V_3$  are  $X$ 's, the remaining  $V$ 's are ringers. This theory happens to have been more or less correct, prior to 1989; see Eq. (42) in Section 12.5.

**THEORY 2.** Again, I use the covariance matrix for the seven observable subtests  $V_1 = \text{NO}, \dots, V_7 = \text{GS}$ , together with the other three unobservable subtests CS, AS, PC. I create an auxiliary variable  $U$ , which is independent of the 10 subtests and has small variance. The distribution of these 11 variables is defined to be jointly normal, and all variables have mean 0. There are three additional unobservables, defined as

$$T_1 = .25(\text{AR} + \text{NO}) + .5\text{PC} + U, \quad (38)$$

$$T_2 = .25(\text{WK} + \text{NO}) + .5\text{PC} + U, \quad (39)$$

$$T_3 = .75(\text{AR} + \text{WK}) - 2U. \quad (40)$$

Let

$$Y = T_1 + T_2 + T_3. \quad (41)$$

In Theory 2,  $T_1, T_2, T_3$  are the unobservables; all the  $V$ 's are ringers. The auxiliary variables  $U$ , CS, AS, PC serve only to define the joint distribution.

Theory 1 and Theory 2 provide the same joint distribution for the observables. Therefore, no statistical procedure based on the joint distribution—like the SGS algorithms or any other correlational methods—can adjudicate between the two theories.

This section and the previous one demonstrate the obvious: you cannot infer cause and effect relationships by doing arithmetic on a correlation matrix, because association is not causation. The mathematical development in SGS avoids such problems only by imposing more or less arbitrary conditions (like faithfulness) on unobservable variables, as discussed in Sections 11.2 and 12.1.

In the present section, neither Theory 1 nor Theory 2 fits into the SGS framework;  $Y$  is a deterministic function of the explanatory variables, with no stochastic error term: see Eq. (35). Furthermore, if  $U$  and PC are treated as variables rather than error terms in (38)–(40), the joint distribution in Theory 2 is, presumably, unfaithful to its causal graph. Similar comments apply to the previous section.

## 12.5. Institutional Background on the AFQT

The “Armed Services Vocational Aptitude Battery” (ASVAB) has 10 subtests, including the seven listed in Table IV, Section 9 above. All 10 are shown in Table VIII.

TABLE VIII  
The 10 Subtests in ASVAB

1. Numerical Operations	NO
2. Word Knowledge	WK
3. Arithmetical Reasoning	AR
4. Mathematical Knowledge	MK
5. Electronics Information	EI
6. Mechanical Comprehension	MC
7. General Science	GS
8. Coding Speed	CS
9. Auto & Shop Information	AS
10. Paragraph Comprehension	PC

*Notes.* The first seven were analyzed by SGS. ASVAB Form 17, July 1990.

Until January, 1989 the AFQT was computed as

$$\text{AFQT} = .5 \times \text{NO} + \text{AR} + \text{WK} + \text{PC}. \quad (42)$$

After that date, NO was replaced by MK; a “verbal” score VE was defined as  $\text{VE} = \text{WK} + \text{PC}$ ; and terms were standardized to have mean 0 and variance 1 on some calibration data—the “NORC 1985 sample.” AFQT was redefined as

$$\text{AFQT} = \text{MK}_Z + \text{AR}_Z + 2 \times \text{VE}_Z, \quad (43)$$

where the subscript Z denotes standardization. Throughout the period, raw scores were by Congressional requirement converted to percentiles based on the NORC sample. Presumably, the data used by SGS [62] come from 1988 or before, since they pick up formula (42) rather than (43); see Section 9 above.<sup>32</sup>

### 13. RESPONSES

Formal statistical inference is, by its nature, conditional. If assumptions A, B, C, . . . hold, then H can be tested against the data. However, if A, B, C, . . . remain in doubt, so must inferences about H. Indeed, the statistical calculations may prove to be quite misleading.

Many assumptions are made but only a few are tested. Those made without testing are called “maintained hypotheses.” They are usually

<sup>32</sup> SGS [62] appear to be considering raw scores, and I follow suit. The material in this section was reported by Larry Hanser (personal communication); he refers to Welsh *et al.* [68, p. 5, Table 3] and Eitelberg [14, p. 73].

statistical and often rather technical—linearity, independence, exogeneity, etc. Careful scrutiny of such assumptions would therefore seem to be a critical part of empirical work.

In the social sciences, however, statistical assumptions are rarely made explicit, let alone validated. Questions provoke reactions that cover the gamut from indignation to obscurantism. *We know all that. Nothing is perfect. Linearity has to be a good first approximation. The assumptions are reasonable. The assumptions do not matter. The assumptions are conservative. You cannot prove the assumptions are wrong. The biases will cancel. We can model the biases. We are only doing what everybody else does. Now we use more sophisticated techniques. What would you do? The decision-maker has to be better off with us than without us. We all have mental models; not using a model is still a model.*

With the SGS approach, responses are more subtle but no more empirical. Proponents often seem to take a Bayesian stance; faithfulness is justified on the grounds that the exceptional cases have measure 0 and must therefore be viewed as negligible *a priori*.<sup>33</sup> However, the SGS approach is frequentist not Bayesian; the simulations, being done on finite-state computers, must concentrate in a set of measure 0; and the SGS class of models has measure 0 within larger classes of models. Indeed, from my perspective, the whole class of path models seems rather unlikely—given the intensity of the research effort and the paucity of convincing examples. The assumptions that diagrams are sparse and faithful stretch credibility even further.

Attempts have also been made to justify the faithfulness assumption by appeals to continuity. If a covariance matrix is unfaithful, small changes to parameter values make it faithful. However, the same argument can be turned against correlational methods. For example, if a covariance matrix is faithful to an incomplete graph, small changes to hidden parameters make the graph complete and vitiate the SGS search procedures. Section 12.1 points to another kind of instability in the SGS framework. The continuity defense (like the Bayesian argument) reflects an aesthetic judgment about modeling styles. Taste is no substitute for empirical verification.

The SGS criteria for causality may also be defended as follows—it is unlikely that anything could produce the patterns of intercorrelation identified by SGS, other than causation; thus, correlational methods shift

<sup>33</sup> The “measure” here is the uniform distribution in Euclidean space, e.g., length, area, volume . . . . The SGS argument [62, p. 95] seems to be a variation on Laplace’s “principle of insufficient reason”: see Stigler [64, p. 127].

the burden of argument. Figures 5 and 6 should dispose of this idea. In real examples, the patterns identified by the SGS search algorithms can hardly represent cause-and-effect relationships. The burden would seem to be on the modelers: how can they recommend an algorithm that gives such results?

Proponents of modeling can also be heard to argue that all of us make assumptions about unobservables. However, what is unobservable with one design may become observable with another. And some investigators still deal with unobservables the hard way—by doing the right studies. For example, take Fisher's "constitutional hypothesis": there may be a genetic factor that predisposes you to smoke and to get lung cancer, heart disease, etc.<sup>34</sup> This putative genetic factor is the unobservable common cause for smoking and illness.

The epidemiologists did not deal with the constitutional hypothesis by introducing special assumptions. Instead, they studied the matter empirically, using data from twin studies. For a recent report on the Swedish twin registry, see Floderus *et al.* [16]. On the Finnish twin registry, see Kaprio and Koskenvuo [31]. Data on the Danish twin registry are fragmentary. There are forthcoming data on the U.S. twin registry, which are quite strong [72]. The numbers on lung cancer are suggestive, but still small—this is a rare disease, even among smokers. The data on heart disease and total mortality, however, make the constitutional hypothesis untenable.

### 13.1. *A Comment from Judea Pearl*

Judea Pearl (personal communication) writes that

Correlation-based model-searching schemes produce causal inferences with only limited guarantees. Yet such schemes have potential, if conducted under conditions that screen out accidental independencies while maintaining structural independencies—for example, longitudinal studies under slightly varying conditions. This assumes, of course, that under such varying conditions the parameters of the model will be perturbed, while its structure remains stable. Maintaining such delicate balance under changing conditions may be hard in real-life studies. However, considering the alternative of resorting to controlled, randomized experiments, such longitudinal studies are still an exciting opportunity.

Additionally, any investigator who is searching for a causal model knowing that the parameters might be tied together by some hidden equation, like (17) [Section 11.2], is wasting time (and public funds). Such a model, even if correct, is bound to be useless, because without the assumption of autonomy (i.e., that each parameter can be perturbed without altering the others), the model cannot predict the effect of interventions or other changes . . . .

Also see Pearl [45]; Pearl and Wermuth [50].

<sup>34</sup> See SGS [62, pp. 298–299].

## 14. OTHER LITERATURE

There is an extensive literature on the evaluation of models, going back at least to the Keynes–Tinbergen exchange (Keynes [32, 33]; Tinbergen [66]). Also see [37, 38]. For more recent discussions, with other citations to the literature, see Freedman [18, 19]. Many authors have tried to explain the basis for inferring causation by using regression. See, for example, [52, 53], or [28, 29]. Of enthusiastic views on social-science modeling, there is no shortage; see, for instance, [60], or [2]. For recent discussions of causal modeling, see Humphreys and Freedman [73], Cox and Wermuth [8], or Pearl [74].

## 15. CONCLUSIONS

SGS [63] have not succeeded in clarifying the circumstances under which causal inferences can be drawn from observed associations, nor have they invented a reliable engine for performing this feat. Their algorithms have some technical interest, but will make causal inferences only when causation is assumed in the first place. To be more explicit: If we assume that the arrows in a path diagram represent causation rather than association, and we also assume that the path diagram can be estimated from data, then indeed SGS can infer causation from association. The faithfulness assumption and exact conditional independence will together eliminate certain kinds of confounding. Even so, causality is assumed into the picture at the beginning, not proved at the end. As Nancy Cartwright says, “No causes in, no causes out.”<sup>35</sup>

The larger problem remains. Can quantitative social scientists infer causality by applying statistical technology to correlation matrices? That is not a mathematical question, because the answer turns on the way the world is put together. As I read the record, correlational methods have not delivered the goods. We need to work on measurement, design, theory. Fancier statistics are not likely to help much.

## APPENDIX: REGRESSION AND CONDITIONING

For ease of reference, this appendix presents the usual formulas for computing regressions and conditional covariances. I begin with regression. Suppose  $\xi$  and  $\eta$  are random variables;  $\xi$  may be a row vector. We seek the column vector  $\beta$  of regression coefficients for  $\eta$  on  $\xi$ . Let

<sup>35</sup> Cartwright [5, Chaps. 2, 3]. Also see Pearl and Verma [49].

$C = E\{\xi\xi'\}$  and  $D = E\{\xi'\eta\}$ ; the prime denotes matrix transposition. Assume  $C$  is positive definite. Then

$$\beta = C^{-1}D. \quad (44)$$

Now  $\eta = \xi\beta + u$ , where  $u$  is automatically orthogonal to  $\xi$ . The mean square of  $u$  may be computed as follows:

$$E(u^2) = E(\eta^2) - \beta' C \beta. \quad (45)$$

If  $\xi$  and  $\eta$  have mean 0, then  $C = \text{cov}(\xi)$  and  $D = \text{cov}(\xi, \eta)$ ; also,  $E(u) = 0$ . Likewise, if some component of  $\xi$  is a nonzero constant,  $E(u) = 0$ . If now the variables are jointly normal,  $u$  is independent of  $\xi$ .

I turn to estimation. Recall Eq. (2), repeated here for ease of reference.

$$Y = X\beta + \epsilon. \quad (2)$$

In this equation,  $X$  is the "design matrix," representing the explanatory variables. There is one row for each unit in the study, and one column for each variable. The entry in the  $i$ th row and  $j$ th column represents the  $j$ th variable, as observed on the  $i$ th unit in the study.  $X$  may include a column of ones if there is to be an intercept in the equation.  $Y$  is a column vector representing the dependent variable, whose  $i$ th component represents the value of  $Y$  for the  $i$ th unit in the study.  $\epsilon$  is also a column vector, with one component for each unit in the study, representing the impact on  $Y$  of chance factors unrelated to  $X$ . Typically, there will be many fewer parameters than data points, so  $\beta$  has relatively few components.

The ordinary least squares estimator for  $\beta$  is denoted by a hat and may be computed as

$$\hat{\beta} = (X'X)^{-1}X'Y. \quad (46)$$

The covariance matrix for  $\hat{\beta}$ , conditional on the design matrix, is computed as

$$\text{cov}(\hat{\beta}|X) = (X'X)^{-1} \text{var}(\epsilon_i|X). \quad (47)$$

Of course, (46) is related to (44); this is seen by defining  $(\xi, \eta)$  as a row chosen at random from  $(X, Y)$ .

The "predicted values" and "residuals" are defined as  $\hat{Y} = X\hat{\beta}$  and  $e = Y - \hat{Y}$ . The residuals are automatically orthogonal to  $X$ . The residual sum of squares, minimized by the choice of  $\beta$ , is  $\text{RSS} = \|e\|^2 = \sum_i e_i^2$ . Then  $\text{var}(\epsilon_i|X)$  in (47) may be estimated as  $\text{RSS}/(n-p)$ , where  $n$  is the number of data points and  $p$  is the number of explanatory variables. Variances will be found along the diagonal of the covariance matrix, and the standard error is computed as the square root of the variance. In



deriving these formulas, it is assumed that, given  $X$ , the components of  $\epsilon$  are conditionally independent and identically distributed, with mean 0.

Suppose the model has an intercept. Then  $R^2$  may be defined as  $R^2 = \text{var}\{\hat{Y}\}/\text{var}\{Y\}$ , where, e.g.,

$$\text{var}\{Y\} = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2, \quad \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i.$$

If all variables have mean 0, then  $R^2 = \hat{\beta}' X' X \hat{\beta} / (n \times \text{var}\{Y\})$ .

The usual formula for computing conditional covariances may be presented as follows. Let  $n > 2$ . Suppose  $X_1, X_2, \dots, X_n$  are jointly normal. We seek the conditional covariance of  $X_1$  and  $X_2$ , given  $X_3, X_4, \dots, X_n$ . Let  $\Sigma$  be the covariance matrix of  $X_3, X_4, \dots, X_n$ . Let  $\kappa_1$  be the covariance of  $X_1$  with  $X_3, X_4, \dots, X_n$ ; let  $\kappa_2$  be the covariance of  $X_2$  with  $X_3, X_4, \dots, X_n$ . We view  $\kappa_1$  and  $\kappa_2$  as  $(n-2) \times 1$  column vectors. The conditional covariance is given by

$$\text{cov}(X_1, X_2 | X_3, \dots, X_n) = \text{cov}(X_1, X_2) - \kappa_1' \Sigma^{-1} \kappa_2. \quad (48)$$

The prime denotes matrix transposition. Details on the material in this appendix may be found in standard texts, for instance, [54].

## ACKNOWLEDGMENTS

Many useful comments were made by Dick Berk, John Cairns, Cliff Clogg, Mark Hansen, Larry Hanser, Jerome Horowitz, Paul Humphreys, Ron Lee, Tony Lin, Bill Mason, Vaughn McKim, Judea Pearl, Diana Petitti, Jamie Robins, Tom Rothenberg, Terry Speed, and Steve Turner. Amos Tversky's work on the paper amounted to collaboration. A version of this paper will also appear in a volume of proceedings, edited by Vaughn McKim and Steve Turner, published by Notre Dame Press.

## REFERENCES

1. L. M. Bartels, Instrumental and "quasi-instrumental" variables, *Amer. J. Pol. Sci.* **35** (1991), 777-800.
2. L. M. Bartels and H. E. Brady, The state of quantitative political methodology, in "Political Science: The State of the Discipline II" (Ada W. Finifter, Ed.), Amer. Pol. Sci. Assoc., Washington, DC, 1993.
3. P. M. Blau and O. D. Duncan, "The American Occupational Structure," Wiley, New York, 1967.
4. J. Cairns, "Cancer: Science and Society," Freeman, San Francisco, 1978.
5. N. Cartwright, "Nature's Capacities and Their Measurement," Clarendon Press, Oxford, 1989.

6. C. C. Clogg and A. Haritou, "The Regression Method of Causal Inference and a Dilemma with This Method," Technical report, Department of Sociology, Pennsylvania State University, 1994.
7. J. Cornfield, W. Haenszel, E. C. Hammond, A. M. Lilienfeld, M. B. Shimkin, and E. L. Wynder, Smoking and lung cancer: Recent evidence and a discussion of some questions, *J. Nat. Cancer Inst.* **22** (1959), 173–203.
8. D. R. Cox and N. Wermuth, Linear dependencies represented by chain graphs, *Statist. Sci.* **8** (1993), 204–283 (with discussion).
9. R. Daggett and D. Freedman, Econometrics and the law: A case study in the proof of antitrust damages, in L. LeCam and R. Olshen, eds. "Proceedings, Berkeley Conference in Honor of Jerzy Neyman and Jack Kiefer," Vol. I, pp. 126–175, Wadsworth, Belmont, CA, 1985.
10. J. N. Darroch, S. L. Lauritzen, and T. P. Speed, Markov fields and log-linear interaction models for contingency tables, *Ann. Statist.* **8** (1980), 522–539.
11. A. Desrosières, "La politique des grands nombres," Éditions la Découverte, Paris, 1993.
12. O. D. Duncan, "Introduction to Structural Equation Models," Academic Press, New York, 1975.
13. A. S. C. Ehrenberg and J. A. Bound, Predictability and Prediction, *J. Roy. Statist. Soc. Ser. A* **156**, Part 2 (1993), 167–206.
14. M. J. Eitelberg, "Manpower for Military Occupations," Office of the Assistant Secretary of Defense (Force Management and Personnel), Washington, DC, 1988.
15. R. F. Engle, D. F. Hendry, and J. F. Richard, Exogeneity, *Econometrica* **51** (1983), 277–304.
16. B. Floderus, R. Cederlof, and L. Friberg, Smoking and mortality: A 21-year follow-up based on the Swedish Twin Registry, *Internat. J. Epidemiology* **17** (1988), 332–340.
17. D. Freedman, A note on screening regression equations, *Amer. Statist.* **37** (1983), 152–155.
18. D. Freedman, As others see us: A case study in path analysis, *J. Educ. Statist.* **12** No. 2 (1987), (with discussion, whole issue).
19. D. Freedman, Statistical models and shoe leather, in "Sociological Methodology 1991" (P. Marsden, Ed.), Amer. Sociol. Assoc., Washington, DC, 1991.
20. D. Freedman and D. Lane, "Mathematical Methods in Statistics," Norton, New York, 1981.
21. C. F. Gauss, "Theoria Motus Corporum Coelestium," Perthes and Besser, Hamburg, 1809; reprinted by Dover, New York, 1963.
22. D. Geiger, "Graphoids: A Qualitative Framework for Probabilistic Inference," Ph.D. dissertation, UCLA, Department of Computer Science, 1990.
23. C. Glymour, A review of recent work on the foundations of causal inference, paper presented at the Notre Dame conference, 1993.
24. C. Glymour, R. Scheines, P. Spirtes, and K. Kelly, "Discovering Causal Structure," Academic Press, New York, 1987.
25. M. Hakama, M. Lehtinen, P. Knekt, A. Aromaa, P. Leinikki, A. Miettinen, J. Paavonen, R. Peto, and L. Teppo, Serum antibodies and subsequent cervical neoplasms: A prospective study with 12 years of follow-up, *Amer. J. Epidemiology* **137** (1993), 166–170.
26. J. Hausman, Specification tests in econometrics, *Econometrica* **46** (1978), 1251–1271.
27. S. L. Hofferth and K. A. Moore, Early childbearing and later economic well-being, *Amer. Soc. Rev.* **44** (1979), 784–815.
28. P. Holland, Statistics and causal inference, *J. Amer. Statist. Assoc.* **81** (1986), 945–960.
29. P. Holland, Causal inference, path analysis, and recursive structural equations models, in "Sociological Methodology 1988," (C. Clogg, Ed.), pp. 449–484, Amer. Sociol. Assoc., Washington, DC, 1988.

30. International Agency for Research on Cancer, "Tobacco Smoking," Monographs on the Evaluation of the Carcinogenic Risk of Chemicals to Humans, Vol. 38, IARC, Lyon, France, 1986.
31. J. Kaprio and M. Koskenvuo, Twins, smoking and mortality: A 12-year prospective study of smoking-discordant twin pairs, *Social Sci. Med.* **29** (1989), 1083–1089.
32. J. M. Keynes, Professor Tinbergen's method, *Econ. J.* **49** (1939), 558–570.
33. J. M. Keynes, Comment on Tinbergen's response, *Econ. J.* **50** (1940), 154–156.
34. H. Kiiveri and T. Speed, Structural analysis of multivariate data: A review, in "Sociological Methodology 1982," (S. Leinhardt, Ed.), Jossey Bass, San Francisco, 1982.
35. E. E. Leamer, Vector autoregressions for causal inference, in "Understanding Monetary Regimes," (K. Brunner and A. Meltzer, Eds.); supplement to the *J. Monetary Econ.*, North-Holland, Amsterdam, 1985.
36. A. M. Legendre, "Nouvelles méthodes pour la détermination des orbites des comètes, Courcier, Paris, 1805; reprinted by Dover, New York, 1959.
37. T. C. Liu, Under-identification, structural estimation, and forecasting, *Econometrica* **28** (1960), 855–865.
38. R. E. Lucas Jr., Econometric policy evaluation: A critique, in "The Phillips Curve and Labor Markets," (K. Brunner and A. Meltzer, Eds.), Carnegie–Rochester Conferences on Public Policy, Vol. 1, pp. 19–64, with discussion, supplementary series to the *J. Monetary Econ.*, North-Holland, Amsterdam, 1976.
39. G. S. Maddala, "Introduction to Econometrics," 2nd ed., McGraw-Hill, New York, 1992.
40. C. F. Manski, Identification problems in the social sciences, in "Sociological Methodology 1993," (P. V. Marsden, Ed.), pp. 1–56, Blackwell, Oxford, 1993.
41. P. Meehl, "Clinical versus Statistical Prediction; A Theoretical Analysis and a Review of the Evidence," University of Minnesota Press, Minneapolis, 1954.
42. K. A. Moore and S. L. Hofferth, Factors affecting early family formation: A path model, *Popul. Environ.* **3** (1980), 73–98.
43. J. Pearl, Fusion, propagation and structuring in belief networks, *Artif. Intell.* **29** (1986), 241–288.
44. J. Pearl, "Probabilistic Reasoning in Intelligent Systems," Morgan Kaufmann, San Mateo, CA, 1988.
45. J. Pearl, Comment: Graphical models, causality and intervention, *Statist. Sci.* **8** (1993), 266–273.
46. J. Pearl, "On the Statistical Interpretation of Structural Equations," Technical Report, Computer Science Department, UCLA, 1994a.
47. J. Pearl, "On the Identification of Nonparametric Structural Equations," Technical Report, Computer Science Department, UCLA, 1994b.
48. J. Pearl, D. Geiger, and T. Verma, The logic of influence diagrams, in "Influence Diagrams, Belief Nets and Decision Analysis," (R. M. Oliver and J. Q. Smith, Eds.), pp. 67–87, Wiley, New York, 1989.
49. J. Pearl and T. Verma, A theory of inferred causation, in "Principles of Knowledge Representation and Reasoning: Proceedings of the Second International Conference (J. A. Allen, R. Fikes, and E. Sandewall, Eds.), pp. 441–452, Morgan Kaufmann, San Mateo, CA, 1991.
50. J. Pearl and N. Wermuth, When can association graphs admit a causal explanation? in "Proceedings, Fourth International Workshop on Artificial Intelligence and Statistics, 1993," pp. 141–150; in "Artificial Intelligence and Statistics" (P. Cheeseman and W. Oldford, Eds.), Springer-Verlag, Berlin, 1994.
51. R. Peto and H. zur Hausen (Eds.), "Viral Etiology of Cervical Cancer," Cold Spring Harbor Laboratory, Banbury Report No. 21, 1986.

52. J. Pratt and R. Schlaifer, On the nature and discovery of structure, *J. Amer. Statist. Assoc.* **79** (1984), 9–21.
53. J. Pratt and R. Schlaifer, On the interpretation and observation of laws, *J. Econ.* **39** (1988), 23–52.
54. C. R. Rao, "Linear Statistical Inference and Its Applications," 2nd ed., Wiley, New York, 1973.
55. R. R. Rindfuss, L. Bumpass, and C. St. John, Education and fertility: Implications for the roles women occupy, *Amer. Sociol. Rev.* **45** (1980), 431–447.
56. R. R. Rindfuss, L. Bumpass, and C. St. John, Education and the timing of motherhood: Disentangling causation, *J. Marriage Family* **46** (1984), 981–984.
57. E. Seneta, Discussion, *J. Educ. Statist.* **12** (1987), 198–201.
58. K. J. Sherman, J. R. Daling, J. Chu, *et al.* Genital warts, other sexually transmitted diseases, and vulvar cancer, *Epidemiology* **2** (1991), 257–262.
59. H. Simon, The meaning of causal ordering, in "Qualitative and Quantitative Social Research," (R. K. Merton, J. S. Coleman, and P. H. Rossi, Eds.), pp. 65–81, Free Press, New York, 1980.
60. N. J. Smelser and D. R. Gerstein, "Behavioral and Social Science: Fifty Years of Discovery," National Academy Press, Washington, DC, 1986.
61. T. P. Speed and H. T. Kiiveri, Gaussian Markov distributions over finite graphs, *Ann. Statist.* **14** (1986), 138–150.
62. P. Spirtes, C. Glymour, and R. Scheines, "Causation, Prediction and Search," Lecture Notes in Statistics, Vol. 81, Springer-Verlag, New York/Berlin, 1993.
63. P. Spirtes, R. Scheines, C. Glymour, and C. Meek, "TETRAD II," Documentation for Version 2.2, Technical Report, Department of Philosophy, Carnegie Mellon University, Pittsburgh, PA, 1993.
64. S. Stigler, "The History of Statistics," Harvard University Press, Boston, 1986.
65. M. Timberlake and K. Williams, Dependence, political exclusion and government repression: Some cross national evidence, *Amer. Sociol. Rev.* **49** (1984), 141–146.
66. J. Tinbergen, "Reply to Keynes," *Econ. J.* **50** (1940), 141–154.
67. T. Verma and J. Pearl, "Causal Networks: Semantics and Expressiveness," in "Uncertainty in AI 4" (R. Shachter, T. S. Levitt, and L. N. Kanal, Eds.), pp. 69–76, Elsevier Science, Amsterdam, 1990.
68. J. R. Welsh, S. K. Kucinkas, and L. T. Curran, "Armed Services Vocational Battery (ASVAB): Integrative Review of Validity Studies," Air Force Human Resources Laboratory Report AFHRL-TR-90-22, 1990.
69. H. White, A heteroskedasticity-consistent estimator and a direct test for heteroskedasticity, *Econometrica* **48** (1980), 817–838.
70. H. White, Maximum likelihood estimation of misspecified models, *Econometrica* **50** (1982), 1–25.
71. G. U. Yule, An investigation into the causes of changes in pauperism in England, chiefly during the last two intercensal decades, *J. Roy. Statist. Soc.* **62** (1989), 249–295.
72. D. Carmelli and W. F. Page, Twenty-four year mortality in World War II US male veteran twins discordant for cigarette smoking, *International Journal of Epidemiology* **25** (1996), 554–559.
73. P. Humphreys and D. Freedman, The grand leap, *Br. J. Phi. Sci.* **47** (1996), 113–123.
74. J. Pearl, Causal diagrams for empirical research, *Biometrika*, **82** (1995), 669–710 (with discussion).
75. N. Muñoz, F. X. Bosch, K. V. Shah, A. Meheus, (eds.) "The Epidemiology of Human Papillomavirus and Cervical Cancer" International Agency for Research on Cancer, Lyon. Distributed in the U.S.A. by Oxford University Press, 1992.