

A Rooster Crow Does Not Cause the Sun To Rise: Review of *Causality: Models, Reasoning and Inference*, by Judea Pearl. Cambridge, UK: Cambridge University Press, 2000. 384 pp. \$39.95 (hardcover).

This book is about the formal analyses of cause-effect relationships between sets of observed events and/or underlying variables related to them. Through ten chapters and an exquisitely well-written epilogue, the author, a prolific computer science specialist, essentially "demystifies" the concept of causality and explains its mathematical, statistical, as well as philosophical implications. Along with brief background material on probability theory and graph theory, the first chapter presents the basic paradigms and major problems of causal analysis and sets the tone for what follows in the subsequent chapters. The most difficult question, what constitutes evidence of a cause-effect relationship in observed data, is discussed in chapter 2, ending with the conceptualization of the validity of any such relationship observed. Chapters 3 and 4 get into deeper theoretical treatments of prediction of direct and indirect effects of actions and policies based on data in the presence of an incomplete understanding of the existence of a cause-effect relationship. Identifiability of cause-effect relationship is the central theme of these chapters. The implications of the calculus of intervention, thus developed, are discussed in the context of applications to social and health science problems in chapter 5 and 6, where the popular constructs of structural equations and confounding are presented. In contrast to the graph theory treatment of detecting the presence of confounding and of identifying critical variables that control the effect of confounding (discussed in chapter 3), chapter 6 presents the difficulties of defining and controlling confounding when statistical criteria are used. The theories of counterfactuals and structural models are presented in chapter 7, through which more rigorous definitions of the concepts introduced earlier in the book are obtained. These include concepts such as causal models, action, causal effects, causal relevance, error terms, and exogeneity. The last three chapters (8 through 10) constitute applications of counterfactual analysis. They include methods of the developing bounds of causal relationship from data of imperfect experiments using combinations of graphical and counterfactual models (chapter 8), identification and interpretation of probability of causation (chapter 9), and a formal explication of the notion of "actual cause" (chapter 10).

Paging through formal definitions of "Markovian Parents" (of an ordered set of variables, p. 14), theorems on "observational equivalence of directed acyclic graphs" (p. 19), mathematical distinctions of "back door" and "front door" adjustments for controlling confounding bias (pp. 78–83), and the like throughout the book, readers of this journal at first glance may wonder why this text is at all relevant for our discipline. However, from titles such as "Segregation analysis reveals a major gene effect controlling systolic blood pressure and BMI in an Israeli population" (Cheng et al. 1998), "Effects of intragenic variability at 3 polymorphic sites of the apolipoprotein B gene on serum lipids and lipoproteins

in a multiethnic Asian population" (Choong et al. 1999), and "RH blood groups and diabetic disorders: Is there an effect on glycosylated hemoglobin level?" (Gloria-Bottini et al. 2000) of recent publications in *Human Biology*, it is obvious that analysis of cause-effect relationships from observational data is an integral element of our research methods. As the title of this review indicates, the thesis of this book is that the inference of cause-effect relationships is not mathematically equivalent to establishing an association between events or variables, although the latter is always observed where a true causal relationship exists. Thus, even though the text is at places highly technical for general readers of this journal, this book is extremely useful for human biologists. Quantitative aptitude is a necessary requirement to understand the logic of the author, although it should not be construed as a drawback of the author's presentation. In contrast, the rigor of the treatment of the subject allows readers to understand scenarios in which association analyses may provide sufficient confidence in inferring an "actual" cause-effect relationship. As more incisive tools of molecular biology are being invoked in the discipline of human biology, researchers should find theoretical materials in this volume that may be used to draw a cause-effect inference from an otherwise "imperfect" study design.

As an example, although the author discussed the subject in the context of applications in social sciences and economics, the analysis of equivalent Structural Equation Models (SEMs) should be of particular appeal to human biologists. This is so because SEM is the underlying analytical formulation of the path analysis that served the purpose of delineating the relative roles of causal factors (e.g., 'genetic' versus 'environmental') in many genetic epidemiological studies (see, e.g., Rao 1991). In this sense, if we consider a path diagram (Li 1975) as a visual representation of a SEM, the author's summary that "scientists need not abandon SEM altogether; they need only abandon the notion that SEM is a method of *testing* causal models" (p. 149) should be considered as an important take-home message to any human biologist who uses path analysis as the tool of deciphering specific causal relationships from data on familial aggregation of biological phenotypes.

On a similar note, the discussion on the Simpson's paradox (chapter 6), not to be confused with the nationally televised infamous Simpson DNA trial, is also of paramount importance in human biology. Since confounding variables and substructures hidden in collected data are commonplace notions in biological studies from which cause-effect relationships are inferred through an association analysis, the mathematics of confounding, collapsibility, and exchangeability, discussed in the chapter, should be helpful for interpretation of data from many human biology study designs.

Likewise, the notions and distinctions of probabilities of necessity (PN), sufficiency (PS), disablement (PD), enablement (PE), and both necessity and sufficiency (PNS), discussed in chapter 9 (see pp. 286–289), have significant relevance in many biomedical inquiries. For example, they are important for understanding the rationale of concepts such as the susceptibility of a population to a

risk factor and their identifiability requirements in epidemiological and clinical trial study designs (see, e.g., Khoury et al. 1989, Cheng 1997). These concepts also play critical roles in determining the statistical strength of enhanced occurrences of untoward health outcomes in the presence of exposures to environmental insults (see pp. 299–300 for interpretation of data on leukemia deaths in children in southern Utah with high and low exposure to radiation from the fallout of nuclear tests in Nevada [Finkelstein and Levin 1990]).

Thus, in terms of audience, students and researchers of applied probability and applied graph theory will obviously be delighted with the contents of this text, since they will find enough food for thought for wider applications of their discipline. The subjects of applications are of even wider interest than the ones explicitly mentioned in the text. From the examples listed above, it is clear that this book should be included at least as a valuable reference resource to classroom discussions of human biology, human ecology, as well as genetic epidemiology. The only handicap is that a full appreciation of the logic discussed will require a mathematical sophistication not generally available among the students of these disciplines.

The text is very well organized in terms of its topical subjects, and the author has drawn from a wide variety of research materials to discuss different applications. Nonetheless, in my opinion, some recent developments of statistical genetics and genetic epidemiology that have relevance to causal analysis did not get their deserved attention. Applications to health science problems are discussed to some extent in chapter 6; two of Wright's seminal works (Wright 1921, 1923) on path analysis are cited, along with the work of Pearson et al. (1899); and some methodological studies in epidemiology are used (e.g., Gail 1986, Khoury et al. 1989, Greenland et al. 1999, Greenland and Robins 1988, Hauk et al. 1991). However, use of path analysis and SEM in genetic epidemiology for delineating causes of familial aggregation, developed and widely used by Morton, Rao, and colleagues (see, e.g., Cloninger et al. 1983) are not discussed in the light of the theory discussed in chapter 5. Likewise, genetic linkage analysis (see, e.g., Terwilliger and Ott 1994) may also be formulated as a cause-effect relationship analysis. It would have been of interest to see an evaluation of the different methodologies of linkage analysis (e.g., the traditional pedigree method, affected sib-pair method, population and family-based association study methods) in terms of the theory (of identifiability of actual linkage) discussed in this text.

For a well-produced text such as this one, it is difficult to find too many flaws. Nevertheless, one typographical error is worth noting. On page 231, the author revisited the data on smoking and lung cancer as an example of inferring causal effects through counterfactual logic. These data were discussed earlier on pages 83–85, in section 3.3.3, and not in section 3.4.3, as stated on page 231. I am sure that this minor flaw is of no consequence, and unless one is interested in the smoking–lung cancer data (section 3.3.3) and, at the same time, intrigued by symbolic derivation of causal effects (section 3.4.3), this error would not have been noticed at all.

In summary, Judea Pearl deserves a high mark for his lucid presentation; he brought clarity to the otherwise confusingly written subject of the literature. This text should be a valuable reference material for researchers who use structural equation models or their symbolic representations. In the era when bioinformatics and computational biology are making entry into human biology research, this book should also be a valuable companion for human biologists attempting to draw causal inference from their research.

RANAJIT CHAKRABORTY

*Human Genetics Center
School of Public Health
The University of Texas Houston Health Science Center
Houston, Texas*

Literature Cited

- Cheng, L.S.-C., G. Livshits, D. Carnelli et al. 1998. Segregation analysis reveals a major gene effect controlling systolic blood pressure and BMI in an Israeli population. *Hum. Biol.* 70:59–75.
- Cheng, P.W. 1997. From covariation to causation: A causal power theory. *Psychol. Rev.* 104:367–405.
- Choong, M.L., S.K. Sethi, and E.S.C. Koay. 1999. Effects of intragenic variability at 3 polymorphic sites of the apolipoprotein B gene on serum lipids and lipoproteins in a multiethnic Asian population. *Hum. Biol.* 71:381–397.
- Cloninger, C.C., D.C. Rao, J. Rice et al. 1983. A defense of path analysis in genetic epidemiology. *Am. J. Hum. Genet.* 35:733–756.
- Finkelstein, M.O., and B. Levin. 1990. *Statistics for Lawyers*. New York, NY: Springer-Verlag.
- Gail, M.H. 1986. Adjusting for covariates that have the same distribution in exposed and unexposed cohorts. In *Modern Statistical Methods in Chronic Disease Epidemiology*, S.H. Moolgavkar and R.L. Prentice, eds. New York, NY: Wiley, 3–18.
- Gloria-Bottini, F., E. Antonacci, N. Bottini et al. 2000. RH blood groups and diabetic disorders: Is there an effect on glycosylated hemoglobin level? *Hum. Biol.* 72:287–294.
- Greenland, S., J. Pearl, and J.M. Robins. 1999. Causal diagrams for epidemiological research. *Epidemiology* 10:37–48.
- Greenland, S., and J. Robins. 1988. Conceptual problems in the definition and interpretation of attributable fractions. *Am. J. Epidemiol.* 128:1185–1197.
- Hauk, W.W., J.M. Heuhaus, J.D. Kalbfleisch et al. 1991. A consequence of omitted covariates when estimating odds ratios. *J. Clin. Epidemiol.* 44:77–81.
- Khoury, M.J., W. D. Flanders, S. Greenland et al. 1989. On the measurement of susceptibility in epidemiologic studies. *Am. J. Epidemiol.* 129:183–190.
- Li, C.C. 1975. *Path Analysis—A Primer*. Pacific Grove, CA: Boxwood Press.
- Pearson, K., A. Lee and L. Bramley-Moore 1899. Genetic (reproductive) selection: Inheritance of fertility in man. *Phil. Trans. Roy. Soc. A* 73:534–539.
- Rao, D.C. 1991. Statistical considerations in applications of path analysis in genetic epidemiology. In *Handbook of Statistics*, Vol. 8, C. R. Rao, and R. Chakraborty, eds. Amsterdam: Elsevier, 63–80.
- Terwilliger, J.D., and J. Ott. 1994. *Handbook of Human Genetic Linkage*. Baltimore: The Johns Hopkins University Press.
- Wright, S. 1921. Correlation and causation. *J. Agri. Res.* 20:557–585.
- Wright, S. 1923. The theory of path coefficients: A reply to Niles' criticism. *Genetics* 8:239–255.