

## TWO BOOKS ON REGRESSION DIAGNOSTICS

D. A. BELSLEY, E. KUH AND R. E. WELSCH, *Regression Diagnostics: Identifying Influential Data and Sources of Collinearity*. John Wiley and Sons, New York, 1980, xv + 292 pages, \$32.95.

R. D. COOK AND SANFORD WEISBERG, *Residuals and Influence in Regression*. Chapman and Hall, New York and London, 1982, x + 230 pages, \$25.00.

Review by A. C. ATKINSON

*Imperial College, London*

My brief in writing this article was not only to review the two books but also to provide a survey of the field of regression diagnostics, including some mention of unsolved problems; these suggestions for research come at the end of the article. In the major part of the article the methods of regression diagnostics are developed with reference to the material in the two books. But first the material must be put in context.

Diagnostic methods are directed at the building and criticism of statistical models. Traditional statistical procedures, such as model fitting and hypothesis testing, are not adequate for this task. For example, rejection of a model does not, *per se*, give guidance as to a more suitable model. Similarly, failure to reject a model does not guarantee that all important aspects of the fit have been adequately examined. One effect of the almost ubiquitous use of the computer in statistics has been the relative ease with which graphs may be produced and used to test models. At the level of text books on regression, the increase in the use of graphical material can be seen by comparing the books of Brownlee (Second edition, 1965), Draper and Smith (1966) and Box, Hunter and Hunter (1978). This trend may be partly due to the increased impact in our society of television as against the printed word. But one of the achievements of diagnostic regression analysis is to provide a framework within which to produce graphs which illustrate both the effect of individual observations on aspects of the fitted model and also ways in which the model is systematically inadequate.

A second consequence of the computer is that complicated statistical analyses can now be routinely performed by scientists with little statistical training or expertise. A major use of diagnostic methods, especially plots, is to call attention to important features of the data which may have been overlooked. The combination of computer and the diagnostic approach may then serve as a substitute for the insight and guile of an experienced statistician. For this combination to be effective the methods have to be easy to program, so that they can readily be added to, or incorporated in, a statistical package. They must also be cheap to compute and the results must be easily intelligible.

The product of a diagnostic analysis may be the identification of an inadequate

---

Received March 1983; revised July 1983.