on the statistical agency, which has guaranteed the confidentiality of the data, because it must consider policies for a range of unknown risks and uses and then police compliance by users. Establishment of a staff of "gatekeepers" could be useful, but could such a staff built essentially to service academic users be justified in the current tight budget climate?

The authors are correct in suggesting that we need to focus more attention on obtaining the informed consent of respondents. I fully support their suggestion that federal agencies need to conduct pilot studies to assess the effectiveness and respondent understanding of the statements used in data collection. In fact, the Bureau of Labor Statistics currently has underway, with IRS sponsorship, research into respondents' understanding of and reaction to the language used in confidentiality statements. Much more work needs to be done in applying the laboratory techniques that combine the cognitive sciences and survey research to assess these issues.

The focus of the Duncan and Pearson article is on data about individuals. But confidentiality problems with establishment data are much more complex than for those about people. For one thing, there are fewer establishments than there are people. Businesses can much more easily be classified into subgroups, often with a very small number of units in the groups of particular interest. In addition, a good deal of information about business establishments is available in publicly accessible files that can be matched to the federal system's file and then used to help to disclose confidential data. Moreover, the value of such data to Duncan and Pearson's data spy might be much greater than the value of the data collected about individuals—to say nothing of those who wish to use such data in prosecution and enforcement.

The risk of disclosure is also generally greater for establishment data than for data about individuals, and the stakes for the company can be quite high when trade secrets or business practices are involved. On the other hand, some data—for example, the number of employees or the identification of major products—may not be sensitive at all to some firms but of great concern to others. There is no simple formula for determining which items are the most sensitive.

The problems involved in finding methods for improving access to microdata on establishments for research purposes are complex and difficult, but the need to find solutions is becoming increasingly necessary. Academic researchers are becoming more and more interested in the use of longitudinal microdata files on business establishments, and access to such data would clearly improve some of the public policy research. Statistical agencies have only just begun thinking about these issues, however, and much more work needs to be done.

Duncan and Pearson are quite right in pointing out that research interests and computational capabilities have led to new and more varied demand for publicly collected data. They are also quite right in pointing to the slow and somewhat negative responses from the nation's primary statistical agencies. But their suggestions, while useful, do not point the way to a quick and clear solution. We in the statistical system strongly believe that the absolute protection of confidentiality tends to assure the cooperation of respondents in voluntary surveys (and most government surveys are based on voluntary cooperation) and enhances the quality of the responses. It is true, however, that statistical agencies have not done all that they could to find ways to provide researchers with the data they need within the practical and legal constraints under which the agencies operate. The article properly challenges the nation's statistical system to revisit the confidentiality practices now in place. In doing so, it serves a valuable function. But the problems we face are real, they are complex and there is no easy and quick solution to them.

# Rejoinder

## George T. Duncan and Robert W. Pearson

While generously acknowledging the centrality of the themes we identify, the discussants quite rightly point to wider issues that should command our attention in the future. To give structure to these issues, we cast the discussants' insights into a set of nested frames. The outer frame encompasses the functional effectiveness of a government statistical system in a diverse society with democratic aspirations. The middle frame delineates the nature of the data that society collects and main-

tains about itself, its institutions and its citizens. The inner frame contains technical developments that affect disclosure-limitation methods. Linking these separate frames is an especially appropriate emphasis on the importance of responsible experimentation.

## EFFECTIVENESS OF GOVERNMENT STATISTICS IN A DIVERSE SOCIETY

In drawing attention to the outer frame of an effective statistical system, Cox emphasizes that confidentiality and data access issues are critically important if government statistical agencies are to serve the diverse components of a society aspiring to democratic principles. One of the most important elements underlying an effective statistical system is public trust. How to create and sustain that trust is, as Cox suggests, a decisive question. Trusting a coin requires examining both sides, head and tail. Just so, trusting a statistical agency requires examining both sides, data collection and data dissemination.

On the data collection side, as Keller-McNulty and Norwood suggest in their remarks, ethical issues are of paramount importance. Rightly, in our view, Norwood stresses the promise of a cognitive laboratory in obtaining a better understanding of respondents' views. Particularly on informed consent questions, the empirical studies of a cognitive laboratory can ascertain both ethical parameters and pragmatic substance about nonresponse and quality of response. In itself, a signed informed consent statement cannot solve the data confidentiality problem, as Keller-McNulty cautions. It helps do so only with confidence that "informed" means informed, and that confidence may in certain cases be built on cognitive laboratory work.

On the data dissemination side, as Norwood notes, research interests and computational capabilities have led to new and more varied demand for data. This demand can only be expected to grow in the future, and not only from the academic community. There will be increased demand for data from Congress, government policy shops, political action groups, business organizations and the media.

We agree that it is not sufficient for statistical agencies to fulfill these two sides of their task well; they must also communicate the importance of reliable and varied statistics to the society, while instilling confidence that they are good stewards of the rights of individual respondents. As Cox suggests, to accomplish this goal, agencies should "go public," engaging public interest groups, advocates for privacy and civil liberties, groups representing the disadvantaged and the press. This strategy may require a basic change in the internal culture of agencies, which Norwood hints may have been too inward in their perspective, to work more closely with their respondent and user constituencies.

Ethically and practically, we agree that users must assume greater responsibility for proper use of data that are provided to them. As Norwood notes, whatever mechanism are developed to ensure this, the practical problems for statistical agencies to "police" users must be addressed. Presumably this would require an appropriate mix of administrative policies, ethical codes within the user communities and legislation.

Norwood argues that legislation has not solved the dilemma, and in some cases laws are contradictory and difficult to apply in particular cases. Cox argues for a maintained focus on legislative issues, particularly in obtaining legal responsibilities of the data user. We take a middle ground: Although ambiguity and contradiction are likely to remain in an area in which an ever-changing balance of rights and interests is sought, specific components of legislation can be revised to reflect a growing, if temporary, consensus concerning pieces of this larger issue.

## NATURE OF THE DATA COLLECTED AND MAINTAINED

Norwood emphasizes the difficulty of confidentiality issues for establishment/institutional data, which we did not address. The problem is complicated because (1) there are fewer establishments than people, (2) lots of publicly available collateral data can be used for matching, (3) the value of attributing information is much greater to the data spy and (4) there are prosecution and enforcement motivations for identification. These complications mean that the risk of disclosure is greater than for surveys of individuals or households and that the stakes to the establishment are higher. In spite of these difficulties, there is special value in longitudinal microdata files on establishments for academic and public policy research.

McNulty addresses the prediction that data will be increasingly maintained in and accessed from computer databases. Statisticians in the federal government need to stay apprised of developments by computer scientists working on data security issues in the private sector.

## STATISTICAL DISCLOSURE LIMITATION

At present, various statistical disclosure limitation methods are used by federal statistical agencies. Some of these apply to aggregate data and some apply to microdata. For aggregate data, both

the Economic Research Service and the National Agricultural Statistics Service of the Department of Agriculture, for example, apply the $(n, k)$ concentration rule, in that published data cells must have at least three observations and no one respondent may represent more than 60% of the total. For microdata, the Bureau of the Census, for example, limits disclosure by use of sample files, not releasing geographic identifiers for areas with small population, not releasing extreme values for continuous variables and not releasing any information that is matchable to administrative record files.

Norwood asks how one measures disclosure risk. Surely, the answer lies outside our article. Fortunately, others have begun work on these issues. Lambert draws a distinction between disclosure risk (the extent to which a released record can be linked to a respondent) and disclosure harm (based on the information revealed when a respondent is linked to a released record). She then develops measures of disclosure risk and of disclosure harm. Some relevant empirical results have been recently obtained by Blien, Wirth and Muller (1990).

Similarly, Keller-McNulty identifies relationship disclosure as an important problem. A clear example of this is the audit rules that the IRS employs. And here again, work is underway or recently completely, for example, by Palley and Simonoff (1986).

Cox examines some of the characteristics of matrix masking, as we have presented it. He finds them a good characterization. Then he extends their generality to data perturbation, multiplication by replacing $X$ by log $X$, categorical data and truncation. He suggests the possibility of a calculus of matrix masks. Masking presents difficult practical problems, as Norwood notes, however, and should be used sparingly and in situations where no other means can provide a desirable level of access to the data. Masking provides no magic answers.

## EXPERIMENTATION

Cox emphasizes the value of experiments in this area, a suggestion that we fully endorse. Certainly we have much to learn, and experience is one of our most reliable guides. This enthusiasm, of course, must be tempered by our commitment to the ethical conduct of experiments and the necessity of not

interrupting normal operating procedures. Perhaps the basic concepts behind Box and Draper's (1969) Evolutionary Operations (EVOP) plans provide some guidance: Take a long-term view of continually improving operations by systematically exploring the effects of changes in controllable variables.

## CONCLUSION

Norwood reminds us that there can be no quick and easy solution to the confidentiality and data access problem. As Isaiah Berlin says in *The Crooked Timber of Humanity*,

> The best that can be done, as a general rule, is to maintain a precarious equilibrium that will prevent the occurrence of desperate situations of intolerable choices—that is the first requirement of a decent society.

But let's take Berlin's formulation as only the first step toward a more optimistic future. Let's venture to nudge things just a bit—sensibly, and in a controlled fashion—so that we might move from precarious equilibrium to precarious equilibrium, each further from the desperate situations of intolerable choices and closer to the hopeful situations of tolerable ones. That we dare this is the second requirement of a decent society.

## ADDITIONAL REFERENCES

ADAM, N. R. and WORTMAN, J. C. (1989). Security-control methods for statistical databases: A comparative study. *ACM Computing Surveys* **21** 515–556.

BLIEN, U., WIRTH, H. and MULLER, M. (1990). Identification risk for microdata stemming from official statistics. International Symposium on Statistical Disclosure Avoidance, Netherlands Central Bureau of Statistics, Voorburg, December 13.

BOX, G. E. P. and DRAPER, N. R. (1969). *Evolutionary Operation*. Wiley, New York.

KELLER-McNULTY, S. and UNGER, E. A. (1991). Database systems: Inferential security. Paper commissioned by National Research Council's Committee on National Statistics Panel on Confidentiality and Data Access. Presented at Conference on Disclosure Limitation Approaches and Data Access, Washington D.C., March 1–2.

LAMBERT, D. (1991). Measures of disclosure risk and harm. Paper commissioned by National Research Council's Committee on National Statistics Panel on Confidentiality and Data Access. Presented at Conference on Disclosure Limitation Approches and Data Access, Washington, D.C., March 1–2.