# Comment

## C. L. Mallows and D. Pregibon

How refreshing to see a paper such as this! The real world is much more complicated and interesting than the artificial one in the textbooks. We are sure that everyone engaged in serious statistical consulting has his or her own set of horror stories. It used to be that statisticians were exhorted to learn their trade by apprenticeship; fortunately there seem to be some general principles, now being identified, that suggest how a consultant may reasonably proceed. Dr. Chatfield has been a leader in identifying these principles. We must hope that teachers and the authors of the next generation of textbooks will give these matters the attention they deserve.

We would like to mention some related work that bears on the subject. Chatfield argues that *strategy*, a largely neglected topic, is at least as important as *techniques* of data analysis. Indeed, as vocal proponents of SPC, the statistical community owns a process that is not in control, the *process* of data analysis. Sometimes it works and produces good quality analyses, and sometimes it does not. Unfortunately we do not understand how to ensure good quality. Part of the problem is that there does not exist a "theory of data analysis" on which to base measurements of quality. Mallows and Walley (1980) and Mallows and Pregibon (1987) attempt to get at the underlying principles of data analysis, with limited success. These articles discuss statistical concepts without appealing to probability models that, in much applied work, are wholly contrived.

A fundamental notion that we have been struggling with concerns "judgments of exchangeability." In its most basic form, this concerns knowing when data may be aggregated and when data may be ignored. Aggregation, even as simple as averaging, requires that the units being aggregated are similar or "exchangeable"; if they are not, aggregation is either meaningless or misleading. In collaboration with Draper and Hodges of RAND, we believe that we have made progress in elucidating the components of exchangeability judgments

in data analysis (Draper, Hodges, Mallows and Pregibon, 1991).

Another reason that the process of data analysis is ill-understood is that it is difficult to think in general terms; every problem seems to have special aspects. Or is it that many scholars do not get involved in real-world (messy) problems, so that thinking and writing about strategy is completely alien? Chatfield provides a few references to discussions of strategic issues, and to them we add Daniel and Wood (1971) for regression (note that a flowchart appears on the page facing page 1!) and Nair and Pregibon (1986) for quality improvement experiments.

Another bibliographical note concerns Chatfield's general guidelines in Section 3. In his wonderful book, Polya (1957) describes the qualitative steps in the mathematical problem solving process. These are (1) understanding the problem; (2) devising a plan; (3) carrying out the plan; and (4) looking back. These four steps translate roughly into Chatfield's guidelines (subsection headings in Section 3). What this indicates to us is that solving problems using statistical methods is not very different from problem solving in general. Thus there is commonality to be exploited and articles such as Chatfield's contribute to understanding the problem solving process.

Our final comments concern Chatfield's remarks on software and its role in data analysis. It almost appears that he is blaming the computer for misleading scientific investigators, who are otherwise thoughtful and thorough. Our view is that software can help rather than hinder the data analysis process. This can happen and is happening in two distinct and complementary ways.

The first focuses on techniques and, in particular, interactive graphical techniques. Indeed, in the quotation attributed to Cox, surely graphical methods are better suited to draw our attention to nonstandard features than are nongraphical methods. Tukey has said "numerical summaries focus on expected values, graphical summaries on unexpected values." The challenge is to harness the computer to work for us. A specific exmaple might help to highlight the issue: In Chatfield's first example concerning "perfectly correlated" variables and the detective work required to explain the anomaly, we envisage a simple interface such that the entries of a correlation matrix are "mouse

*C. L. Mallows is a member of the technical staff and D. Pregibon is Department Head at AT&T Bell Laboratories, 600 Mountain Avenue, Murray Hill, New Jersey 07974-2070.*