

In This Issue

The following type of example is often presented in introductory probability and statistics courses to help sharpen students' intuition about the importance of background rates in calculating probabilities: Suppose that you are walking down the street and notice that the Department of Public Health is giving a free medical test for a certain rare disease. The test is 90% reliable in the following sense: If a person has the disease, there is a probability of 0.9 that the test will give a positive response (the "sensitivity" of the test); and if a person does not have the disease, there is a probability of 0.9 that the test will give a negative response (the "specificity" of the test). Data indicate that your chances of having the disease are only 1 in 5000. However, because the test costs you nothing (you have already paid for it with your taxes), and it is fast and harmless, you decide to stop and take the test. A few days later you learn that you had a positive response to the test. What is now the probability that you have the disease?

Many beginning students feel that this probability should be about 0.9, but that feeling mistakenly ignores the small prior probability of 0.0002 that you had the disease. The correct posterior probability is found by Bayes theorem to be 0.0018. Your probability of having the disease is now 9 times as large as it was before you took the test, but it is still extremely small. The intuitive explanation is that because the test has a 10% rate of producing false positives, there will be about 500 positive responses among a group of 5000 persons, but on the average only one person in the group will have the disease.

It is this large number of false positives that has led various interested parties to question the effectiveness of large-scale medical screening tests for populations in which the prevalence of the disease is low, and which is the subject of the opening article by Joseph L. Gastwirth in this issue. He considers problems in which the prevalence, as well as the sensitivity and the specificity of the test, are unknown, and discusses effective experimental designs for estimating these quantities in order to obtain an estimate of the posterior probability given a positive response that will have small variance. He describes two applications that have been very much in the news in recent years: the screening of general populations for the presence of antibodies to the AIDS virus and the screening of the employees, or potential employees, of an organization with polygraph (or "lie detector") tests.

In his discussion of this article, D. H. Kaye considers the standards that are used for the admissibility of polygraph evidence in court, and the relevance of

Gastwirth's work to the legal question of admissibility. John C. Kircher and David C. Raskin point out that the problem of low base rates has been discussed for many years in the psychology literature, and describe the many different contexts in which polygraph tests are used. Janet Wittes emphasizes that the context of a medical screening determines whether the sensitivity or the specificity of the test is more important. Judith D. Goldberg points out that not only prevalence, but also false positive and false negative rates, can vary from group to group. Seymour Geisser sketches a Bayesian predictive approach to the problems addressed by Gastwirth. Finally, Beth C. Gladen comments that in many situations the application of a confirmatory test following a positive response would make variance calculations relatively unimportant.

* * *

In his article, "Uncertainty, policy analysis, and statistics," James S. Hodges states that "No existing school of statistical thinking provides a comprehensive framework for considering the various types of uncertainty and the tradeoffs among them that analysts must make." He describes three major types of uncertainty: (1) structural uncertainty, which is uncertainty about the model that is used; (2) risk, which is uncertainty due to statistical or stochastic variability given the model; and (3) technical uncertainty, which is uncertainty due to data processing and the use of approximations. He argues that the absence of a system that properly accounts for all these types "creates an inherent tendency for analyses to understate uncertainty about predictions . . . which can lead to invisible biases in policy considerations." He believes that the de Finetti approach comes closest to providing such a system, and he tries in this paper to develop further the connection between that approach and real policy applications.

In his comment, David Freedman states that "Good statistical analysis can be done in either the frequentist or the Bayesian framework. However, for either approach to succeed, the analyst has to get the model right, or close enough." Seymour Geisser points out that there is a fundamental principle in the de Finetti approach to statistics "that statisticians (even Bayesian predictivists) often ignore." Peter J. Huber comments on "the problem of the infinite regress, and the question of whether and when to combine different kinds of uncertainty." Joseph B. Kadane stresses three important aspects of de Finetti's approach: the subjectivity of probability, the emphasis on prevision and

the insistence on finitely additive probabilities. Albert Madansky uses "more earthy terms" to describe the three major types of uncertainty covered by Hodges: "The model . . . may be 'off-base,' the procedures recommended . . . may be 'dead wrong,' and . . . the statistician may 'drop the ball' in implementing his recommended procedure." He then goes on to raise (and answer) some questions suggested by the three terms in the title of Hodges' paper. Adrian F. M. Smith points out that "Model development and use typically involves a progression through five broad stages: perceived problem situation, conceptual model, formal model, technical solution and summary output of some kind." He emphasizes that "the quantitative analyst cannot and should not be acting in splendid technical isolation" but instead must be part of a team with broad expertise.

* * *

In many empirical investigations to learn about the effects of a new treatment, it is not possible or not feasible for the investigator to decide, either randomly or deterministically, in the course of the experiment which subjects will receive the treatment and which ones will not. Instead, the subjects that might be included in the study can only be classified as either having received the treatment or not. An *observational study* is an attempt to estimate the effects of the treatment in an investigation of this type. In his article in this issue, Paul R. Rosenbaum studies the usefulness of more than one control group in observational studies. He writes, "A second control group provides a test of the assumption that conventional adjustments for observed covariates suffice in estimating treatment effects."

In his discussion, Paul W. Holland describes his experience with multiple control groups in an observational study of computer-assisted instruction. Barry H. Margolin discusses the use of multiple control groups in designed or controlled experimentation. Richard G. Cornell extends the discussion of case-control studies and observational studies comparing groups having different exposures to harmful agents. Norman Breslow states that "The major conclusions of this paper should come as no surprise to biostatisticians and epidemiologists involved in the applications of statistical methods and concepts to clinical and observational studies in public health and medicine" because the principles of study design and interpretation "are well known and widely used." Rosenbaum, in his rejoinder, disagrees.

* * *

In problems of testing hypotheses about a parameter θ , a precise hypothesis is one which specifies that θ has a particular value or that θ lies in a given small

interval. In their article, James Berger and Mohan Delampady review the problem of testing a precise null hypothesis, "with special emphasis placed on exploring the dramatic conflict between conditional measures (Bayes factors and posterior probabilities) and the classical P-value (or observed significance level)." Their studies show that "claiming that a P-value of 0.05 is significant evidence against a precise hypothesis is sheer folly; the actual Bayes factor may well be near 1, and the posterior probability of H_0 near $\frac{1}{2}$." This leads them to the recommendation that "when testing precise hypotheses, formal use of P-values should be abandoned. Almost anything will give a better indication of the evidence provided by the data against H_0 ." The paper is clearly and enjoyably written, and with its strong conclusions it has inevitably sparked interesting discussion.

D. R. Cox writes that "the paper is a valuable and thought-provoking one," but "the conclusion that P-values have no role at all is wrong." He also gives his own maxim to accompany others in the article: "Attempts to force formal problems of statistical inference into an exclusively Bayesian mold may give misleading answers." Morris L. Eaton discusses whether there is a specific question that a P-value addresses; the use of "automatic" methods, Bayesian or frequentist; and the meaning of objectivity. Arnold Zellner suggests that in many problems, prior information distinguishing certain important alternatives is usually available and that Bayes factors should be computed for them "to get a hold on the sensitivity of results to specific, relevant broader assumptions." M. J. Bayarri presents a Bayesian decision-theoretic approach to goodness-of-fit tests in which the special nature of the precise null hypothesis is represented in the utility function, because of the usefulness of the particular model being tested, rather than in a sharply spiked prior distribution. George Casella and Roger L. Berger disagree with the authors and argue that "it is not the case that P-values are too small, but rather that Bayes point null posterior probabilities are much too big!" Joseph B. Kadane notes that a significance test is driven by the sample size and concludes "that the technique of testing hypotheses is vastly overrated in statistics as a method."

Readers who crave still more discussion on this topic are referred to the March 1987 issue of the *Journal of the American Statistical Association*, pages 106-139.

* * *

The featured interview in this issue of *Statistical Science* is a conversation with George E. P. Box, Vilas Research Professor of Mathematics and Statistics at the University of Wisconsin, who is well-known for his work in experimental design, time series analysis, quality improvement and other areas of statistics.