

Why Isn't Everyone a Bayesian?

B. EFRON*

Originally a talk delivered at a conference on Bayesian statistics, this article attempts to answer the following question: why is most scientific data analysis carried out in a non-Bayesian framework? The argument consists mainly of some practical examples of data analysis, in which the Bayesian approach is difficult but Fisherian/frequentist solutions are relatively easy. There is a brief discussion of objectivity in statistical analyses and of the difficulties of achieving objectivity within a Bayesian framework. The article ends with a list of practical advantages of Fisherian/frequentist methods, which so far seem to have outweighed the philosophical superiority of Bayesianism.

KEY WORDS: Fisherian inference; Frequentist theory; Neyman-Pearson-Wald; Objectivity.

1. INTRODUCTION

The title is a reasonable question to ask on at least two counts. First of all, everyone used to be a Bayesian. Laplace wholeheartedly endorsed Bayes's formulation of the inference problem, and most 19th-century scientists followed suit. This included Gauss, whose statistical work is usually presented in frequentist terms.

A second and more important point is the cogency of the Bayesian argument. Modern statisticians, following the lead of Savage and de Finetti, have advanced powerful theoretical reasons for preferring Bayesian inference. A byproduct of this work is a disturbing catalogue of inconsistencies in the frequentist point of view.

Nevertheless, everyone is not a Bayesian. The current era is the first century in which statistics has been widely used for scientific reporting, and in fact, 20th-century statistics is mainly non-Bayesian. [Lindley (1975) predicts a change for the 21st!] What has happened?

2. TWO POWERFUL COMPETITORS

The first and most obvious fact is the arrival on the scene of two powerful competitors: Fisherian theory and what Jack Kiefer called the Neyman-Pearson-Wald (NPW) school of decision theory, whose constituents are also known as the frequentists. Fisher's theory was invented, and to a remarkable degree completed, by Fisher in the period between 1920 and 1935. NPW began with the famous lemma of 1933, asymptoting in the 1950s, though there have continued to be significant advances such as Stein estimation, empirical Bayes, and robustness theory.

Working together in rather uneasy alliance, Fisher and NPW dominate current theory and practice, with Fisherian ideas particularly prevalent in applied statistics. I am going to try to explain why.

3. FISHERIAN STATISTICS

In its inferential aspects Fisherian statistics lies closer to Bayes than to NPW in one crucial way: the assumption that there is a *correct* inference in any given situation. For example, if x_1, x_2, \dots, x_{20} is a random sample from a Cauchy distribution with unknown center θ ,

$$f_{\theta}(x_i) = \frac{1}{\pi[1 + (x_i - \theta)^2]}$$

then in the absence of prior knowledge about θ the correct 95% central confidence interval for θ is, to a good approximation,

$$\hat{\theta} \pm 1.96 / \sqrt{-\ddot{l}_{\hat{\theta}}}$$

where $\hat{\theta}$ is the maximum likelihood estimator (MLE) and $\ddot{l}_{\hat{\theta}}$ is the second derivative of the log-likelihood function evaluated at $\theta = \hat{\theta}$. The (mathematically) equally good approximation

$$\hat{\theta} \pm 1.96/\sqrt{I_0}$$

(10 being the expected Fisher information), is not correct (Efron and Hinkley 1978).

Fisher's theory is a theory of archetypes. For any given problem the correct inference is divined by reduction to an archetypal form for which the correct inference is obvious. The first and simplest archetype is that of making inferences about θ from one observation x in the normal model

$$x \sim N(\theta, 1). \quad (1)$$

Fisher was incredibly clever at producing such reductions: sufficiency, ancillarity, permutation distributions, and asymptotic optimality theory are among his inventions, all intended to reduce complicated problems to something like (1). (It is worth noting that Fisher's work superseded an earlier archetypal inference system, Karl Pearson's method of moments and families of frequency curves.)

Why is so much of applied statistics carried out in a Fisherian mode? One big reason is the *automatic nature* of Fisher's theory. Maximum likelihood estimation is the original jackknife, in Tukey's sense of a widely applicable and dependable tool. Faced with a new situation, the working statistician can apply maximum likelihood in an automatic fashion, with little chance (in experienced hands) of going far wrong and considerable chance of providing a nearly optimal inference. In short, he does not have to think a lot about the specific situation in order to get on toward its solution.

Bayesian theory requires a great deal of thought about the given situation to apply sensibly. This is seen clearly in the efforts of Novick (1973), Kadane, Dickey, Winkler, Smith, and Peters (1980), and many others to at least partially automate Bayesian inference. All of this thinking is admirable in principle, but not necessarily in day-to-day practice. The same objection applies to some aspects of

*B. Efron is Professor, Department of Statistics, Stanford University, Stanford, CA 94305.

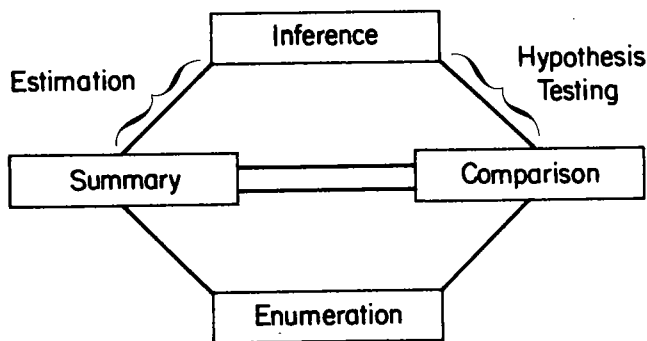


Figure 1. Four Basic Statistical Operations and How They Relate to Estimation. Source: Efron (1982b, fig. 2).

NPW theory, for instance, minimax estimation procedures, and with the same result: they are not used very much.

Not all of statistics is inference. The little diagram of all of statistics in Figure 1 (reprinted from Efron 1982b) starts at the bottom with "enumeration," the collecting and listing of individual datum. The diagram proceeds upward to the reduction of the raw data to more understandable form through the adversarial processes of summary and comparison. This is the part of the analysis where, usually, the statistician decides on a reasonable probabilistic model for the situation. At the top of the diagram is inference. This is the step that takes us from the data actually seen to data that might be expected in the future.

Bayesian theory concentrates on inference, which is the most glamorous part of the statistical world, but not necessarily the most important part. Fisher paid a lot of attention to the earlier steps of the data analysis. Randomization for instance, and experimental design in general, is a statement about how data should be collected, or "enumerated," for best use later in the analysis. Maximum likelihood is a provably efficient way to summarize data, no matter what particular estimation problems are going to be involved in the final inference (Efron 1982b). The NPW school has also contributed to the theory of enumeration, notably in the areas of survey sampling and efficient experimental design.

Fisher's theory culminated in fiducial inference, which to me and most current observers looks like a form of objective (as opposed to subjective) Bayesianism. I will discuss the problems and promise of objective Bayesianism later, but it is interesting to notice that fiducial inference is alone among Fisher's major efforts in its failure to enter common statistical application. In its place, the NPW theory of confidence intervals dominates practice, despite some serious logical problems in its foundations.

4. THE NPW SCHOOL

Unlike Bayes and Fisher, the NPW school does not insist that there is a correct solution for a given inferential situation. Instead, a part of the situation deemed most relevant to the investigator is split off, stated in narrow mathematical fashion, and it is hoped, solved. For example, the correct Bayesian or Fisherian inference for θ in situation (1) leads directly to the correct inference for $\gamma \equiv 1/(1 + \theta)$, but this is not necessarily the case in the NPW formulation. (What is the uniform minimum variance unbiased estimate of γ ?)

The NPW piecewise approach to statistical inference has

been justly criticized by Bayesians as self-contradictory, inconsistent, and incoherent. The work of Savage, de Finetti, and their successors shows that no logically consistent inference maker can behave in such a non-Bayesian way. The reply of the NPW school is that there is no reason to believe that statistical inference should be logically consistent in the sense of the Bayesians, and that there are good practical reasons for approaching specific inference problems on an individual basis.

As an example consider the following problem: we observe a random sample x_1, x_2, \dots, x_{15} from a continuous distribution F on the real line and desire an interval estimate for θ , the median of F . The experiment producing the x_i is a new one, so very little is known about F .

A genuine Bayesian solution seems difficult here, since it requires a prior distribution on the space of all distributions on the real line. Frequentist theory produces a simple solution in terms of a confidence interval based on the order statistics of the sample,

$$\theta \in [x_{(3)}, x_{(12)}]$$

with probability .963, no matter what F may be. The fact that this solution, unlike a Bayesian one, does not also solve the corresponding problem for say $\phi \equiv 50\%$ trimmed mean of F does not dismay the frequentist, particularly if a satisfactory Bayesian solution is not available.

The Bayesian accusation of incoherency of the frequentist cuts both ways: in order to be coherent Bayesians have to solve all problems at once, an often impossible mental exercise.

As another example consider "rejecting at the .05 level." The inconsistencies of this practice are well documented in the Bayesian literature (see Lindley 1982). On the other hand it is one of the most widely used statistical ideas. Its popularity is founded on a crucial practical observation: it is often easier to compare quantities than to assign them absolute values. In this case the comparison is between the amount of evidence against the null hypothesis provided by different possible outcomes of the data. For testing $H_0: x \sim N(0, 1)$ versus $H_1: x \sim N(2, 1)$, we know that a large observed x provides greater evidential value against H_0 and in favor of H_1 , even if we cannot absolutely quantify "evidence."

A Bayes solution to this problem, "the a posteriori odd ratio is 7 to 1 in favor of H_1 ," is more satisfactory than "the data are significant at the .05 level," but it also requires more input. In fact, it tacitly implies that we have assigned an absolute measure of evidence to every possible outcome. Absolute here means that the meaning of 7 to 1 is the same no matter what experiment it came from. [Good's (1965) Bayes-non-Bayes compromise suggests using Bayesian ideas in a comparative mode, but this is the only example I know.]

The heart attack decision tree (Fig. 2) illustrates another difficult situation for the honest Bayesian. The tree purports to predict coronary patients with high risk of dying (population 2) on the basis of variables observed at hospital admission. A series of dichotomous observations are made for example, high or low kinase level, which result in a final prediction. The nodes marked "2" on the tree predict death. Of the 389 patients classified by the tree, only 1 out

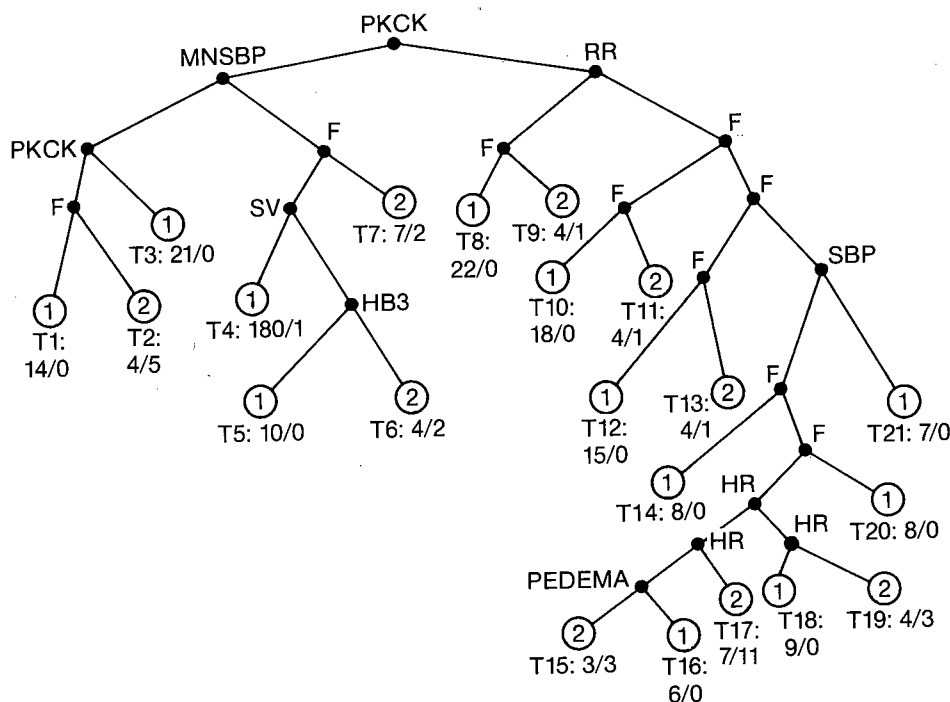


Figure 2. A Decision Tree for Classifying Heart Attack Patients Into Low Risk of Dying (population 1) or High Risk of Dying (population 2). Smaller values of the decision variables go to the left. Circled numbers at terminal nodes indicate population prediction. For example, 6 of the 389 patients in the training set end up at T6, 4 from population 1, 2 from population 2; these patients would all be predicted to be in population 2. Abbreviations: PKCK, peak creatinine kinase level; MNSBP, minimum systolic blood pressure; SBP, systolic blood pressure; FF, respiration rate; HR, average heart rate; SV, supraventricular arrhythmia; HB3, heart block 3rd degree; PEDEMA, peripheral edema; F, Fisher linear discriminant function, differing from node to node. Source: Efron (1982a, fig. 7.1).

of 30 deaths was misclassified, that is, predicted to live. Can we believe that the tree has 96.7% probability of successfully predicting deaths?

Because the medical investigators had little prior knowledge of the situation, the tree was constructed by an elaborate data-fitting procedure, which in fact was designed to maximize the apparent success rate. At each stage the dichotomous variable to be used and the splitting point defining "high" or "low" were chosen to give the maximum apparent difference between populations 1 and 2. A bootstrap analysis, much like a cross-validation, gave an unbiased estimate of successful prediction of death of about 70%, rather than 96.7%, for this tree. (Details appear in sec. 7.6 of Efron 1982a.)

The fact that the observed data were used to construct the tree, and how they were used, makes no difference to the Bayesian, since it has no effect on the likelihood function. This is similar in spirit to the fact that the stopping rule used in a sequential procedure has no Bayesian consequence. It makes a world of difference to the frequentist. If exactly the same tree had been constructed by a less flexible rule, the unbiased estimate would move closer to the observed value 96.7%. This is incoherent behavior. The Bayesian estimate, whatever it is, would not change.

"Ad hoc" is a pejorative adjective in Bayesian descriptions of frequentist statistics. On the other hand, ad hoc reasoning produces a reasonable answer here, in a problem that seems far too complicated for a full Bayesian solution. The right to split off the simple part of a complicated inference problem should not be the exclusive property of the frequentists, but I am not aware of much Bayesian activity

along these lines. The coherency approach of Savage and de Finetti seems to have discouraged it. (For a counterexample to this statement see Boos and Monahan, in press.)

The NPW school invented decision theory, but it is not decision theory that separates them from the Bayesians. In fact, Bayesians have made good use of decision theory. The parting of the ways occurs on the crucial issue of θ averages, expectations taken with the state of nature θ fixed. In other words, frequentist calculations. Controlling, or at least computing, θ averages is central to the NPW approach and irrelevant to the Bayesians. This brings us to the topic of objectivity, in my opinion the linchpin of non-Bayesian success with statistical practitioners.

5. OBJECTIVITY

So far I have been careful not to define the kind of Bayesian theory under criticism. The dominant Bayesian school, and the one with the legitimate claim to philosophic coherency, is the subjective Bayesianism of de Finetti and Savage. Now by definition one cannot argue with a subjectivist, so I will just state the obvious fact: though subjectivism is undoubtedly useful in situations involving personal decision making, for example, business and legal decisions, it has failed to make much of a dent in scientific statistical practice. The nature of scientific communication makes me doubt that it ever will.

"Scientific objectivity" is more than a catch-phrase. Strict objectivity is one of the crucial factors separating scientific thinking from wishful thinking. Complete objectivity about one's own work is a little much to expect from a human being, even a scientist, but it is not too much to expect from

