

Robert E. KASS and Adrian E. RAFTERY*

In a 1935 paper and in his book *Theory of Probability*, Jeffreys developed a methodology for quantifying the evidence in favor of a scientific theory. The centerpiece was a number, now called the *Bayes factor*, which is the posterior odds of the null hypothesis when the prior probability on the null is one-half. Although there has been much discussion of Bayesian hypothesis testing in the context of criticism of *P*-values, less attention has been given to the Bayes factor as a practical tool of applied statistics. In this article we review and discuss the uses of Bayes factors in the context of five scientific applications in genetics, sports, ecology, sociology, and psychology.

We emphasize the following points:

- From Jeffreys' Bayesian viewpoint, the purpose of hypothesis testing is to evaluate the evidence in favor of a scientific theory.
- Bayes factors offer a way of evaluating evidence *in favor of* a null hypothesis.
- Bayes factors provide a way of incorporating external information into the evaluation of evidence about a hypothesis.
- Bayes factors are very general and do not require alternative models to be nested.
- Several techniques are available for computing Bayes factors, including asymptotic approximations that are easy to compute using the output from standard packages that maximize likelihoods.
- In "nonstandard" statistical models that do not satisfy common regularity conditions, it can be technically simpler to calculate Bayes factors than to derive non-Bayesian significance tests.
- The Schwarz criterion (or BIC) gives a rough approximation to the logarithm of the Bayes factor, which is easy to use and does not require evaluation of prior distributions.
- When one is interested in estimation or prediction, Bayes factors may be converted to weights to be attached to various models so that a composite estimate or prediction may be obtained that takes account of structural or model uncertainty.
- Algorithms have been proposed that allow model uncertainty to be taken into account when the class of models initially considered is very large.
- Bayes factors are useful for guiding an evolutionary model-building process.
- It is important, and feasible, to assess the sensitivity of conclusions to the prior distributions used.

KEY WORDS: Bayesian hypothesis tests; BIC; Importance sampling; Laplace method; Markov chain Monte Carlo; Model selection; Monte Carlo integration; Posterior model probabilities; Posterior odds; Quadrature; Schwarz criterion; Sensitivity analysis; Strength of evidence.

1. INTRODUCTION

The Bayesian approach to hypothesis testing was developed by Jeffreys (1935, 1961) as a major part of his program for scientific inference. Although Jeffreys called his methods "significance tests," apparently borrowing the term from Fisher, this is misleading, because Jeffreys's perspective and goals were quite different. Jeffreys was concerned with the comparison of predictions made by two competing scientific theories. In his approach, statistical models are introduced to represent the probability of the data according to each of the two theories, and Bayes's theorem is used to compute the posterior probability that one of the theories is correct.

Considerable attention has been given to distinctions between the two approaches (e.g., Berger and Delampady 1987, Berger and Berry 1988, and references therein). Often lost from the controversy, however, are the practical aspects of

the Bayesian methods: how conclusions may be drawn from them, how they can provide answers when non-Bayesian methods are hard to construct, what their strengths and limitations are. These concerns are the focus of this article. We will also discuss the Bayesian approach to accounting for uncertainty in the model-building process, which is closely connected to the methodology for hypothesis testing.

In Section 2 we motivate the work with several applications from the areas of genetics, sports, ecology, sociology, and psychology. These help connect hypothesis testing with model selection and introduce several problems that Bayesian methodology can solve, including the evaluation of the evidence in favor of a null hypothesis, the inclusion of other information in the weighing of evidence, the comparison of nonnested models, and accounting for uncertainty in the choice of models. In Section 3 we introduce the *Bayes factor*, which is the posterior odds of one hypothesis when the prior probabilities of the two hypotheses are equal.

Bayesian methods involve integrals and thus, often, numerical integration. Many integration techniques have been adapted to problems of Bayesian inference, including the computation of Bayes factors; this is discussed in Section 4. Bayes factors require priors on the parameters appearing in the models that represent the competing hypotheses. The choice of these priors and the extent to which Bayes factors are sensitive to this choice is discussed in Section 5.

* Robert E. Kass is Professor, Department of Statistics, Carnegie Mellon University, Pittsburgh, PA 15213. Adrian E. Raftery is Professor of Statistics and Sociology, Department of Statistics, University of Washington, Seattle, WA 98195. Kass's research was supported by National Science Foundation Grant DMS-9005858 and by National Institutes of Health Grant RO1-CA54852-01. Raftery's research was supported by Office of Naval Research Contract N-00014-91-J-1074, by the Ministère de la Recherche et de l'Espace, Paris, by the Université de Paris VI, and by INRIA, Rocquencourt, France. Raftery thanks the latter two institutions, Paul Deheuvels, and Gilles Celeux for hearty hospitality during his Paris sabbatical in which part of this article was written. The authors are grateful to former editor Don Guthrie for encouraging them to write this article, to David Madigan and Larry Wasserman for many helpful comments and discussions, and to Jim Albert, James Dickey, Andrew Gelman, Julia Mortera, Michael Newton, Sue Rosenkranz, Michael Sobel, Mike Titterington, the editor, the associate editor, and two anonymous referees for very helpful comments on an earlier version of the article.

In Section 6 we take up the problem of accounting for uncertainty about model form. The data analyst is often faced with many models that involve different assumptions, distributional forms, or sets of covariates. Although he or she may wish to summarize findings with a single model, there are usually many choices to be made, and in estimating quantities of interest it is desirable to provide an assessment of the uncertainty that accounts for the model-building process itself. This can be done with Bayesian methods, in which Bayes factors are used to calculate the posterior probabilities of the models considered.

In Section 7 we return to the applications and show how the methods reviewed in Sections 3–6 may be used to solve the problems posed in Section 2.

There has been much controversy about Bayes factors. In Section 8 we discuss several issues, including the purpose of testing sharp hypotheses, disagreements between Bayes factors and P values, and alternative model selection criteria. In Section 9 we briefly mention some other work and provide references to the calculation of Bayes factors for some specific models. Finally, in Section 10 we conclude by summarizing the most important points and highlighting outstanding problems.

2. APPLICATIONS

In this section we present five applications that pose problems usefully solved with Bayes factors. In the first two the goal is to evaluate the evidence *in favor of* a null hypothesis. The third involves irregular, nonnested models. The fourth has to do with drawing inferences while trying to account for the uncertainty in modeling; the fifth, with determination of which of two sets of alternative explanatory variables better predicts some binary repeated-measures data. The latter will lead to computational difficulties solved by some of the techniques reviewed in this article. Here we describe the problems; in Section 7 we describe solutions to them.

Application 1: *Escherichia coli* Mutagenesis

In an experiment in molecular biology (Sklar and Strauss 1980), the investigators hypothesized that in the *uvrE* strain of *E. coli* bacteria, mutations leading to “acetate utilization deficiency,” would occur by an unusual error-prone DNA repair mechanism. As a consequence, this mutation would fail to be linked to mutations at neighboring loci. Specifically, they noted that if the acetate utilization deficiency mutation occurred during DNA replication, then it would be linked to the relatively rare trait of “rifampin resistance,” but if the error-prone repair mechanism were responsible, then there would be no such linkage. The investigators thus created a pair of cell lines, of which one contained cells “selected” for rifampin resistance, and the other contained “unselected” cells. The absence of linkage, predicted by the error-prone repair hypothesis, would imply that the proportions p_1 and p_2 of bacteria exhibiting acetate utilization deficiency in the “selected” and “unselected” cell lines would be equal. When the investigators took samples from each cell line and found \hat{p}_1 and \hat{p}_2 to be approximately equal, they believed that this ought to have represented fairly strong evidence in favor of error-prone DNA repair.

Already there is an interesting problem here. The investigators understood that the hypothesis of no linkage corresponded to the statistical null hypothesis $H_0: p_1 = p_2$, and they computed Pearson’s chi-squared statistic to test it. They found that the chi-squared test did not reject the null hypothesis $p_1 = p_2$, but they were aware that the usual interpretation of significance tests is that they may be used only to *reject* hypotheses and do not offer an assessment of the strength of the evidence *in favor of* the null hypothesis. Thus they were left in doubt about the question that they thought their data should have been able to answer.

In fact, there is more to the story. They also had data on 12 other strains of *E. coli*, which showed a range of variation in the differences between p_1 and p_2 . There clearly was considerable related information provided by the other strains, and it was desirable to try to use that information in the statistical analysis of the *uvrE* data.

This first application poses several problems. First, the investigators wished to know the *strength of the evidence* provided by the data in favor of their scientific hypothesis, which was translated to a statistical hypothesis. Second, the hypothesis at issue was the *null*, and the failure of the Pearson chi-squared statistic to reject it did not indicate the strength of evidence in its favor. Third, there were other data that could be used to construct an alternative hypothesis; that is, there was prior information. These facets of the problem suggest that the Bayesian approach could be useful.

In the next application, evidence in favor of the null is again at issue, but prior information is less well documented. This example concerns the evidence against extrabinomial variability.

Application 2: The Hot Hand

It has been argued that belief in the “hot hand” in basketball is based on misperception of random sequences (Gilovich, Vallone, and Tversky 1985). That is, erratic behavior of shooting may not reflect any real tendency for players to have good streaks or bad streaks but may instead be consistent with a stable shooting percentage. If there were a strong tendency for players to have good and bad streaks, then one would expect to see this expressed in good and bad days. In fact, many people who engage in athletic activity do seem to think that they themselves have good days and bad days. Fans often have similar beliefs about players. As a check on results reported by Gilovich et al. (1985), one of us (Kass) collected data on Larry Bird’s game-by-game performance during the 1986–1987 season. Bird’s field goal shooting percentage ranged from 21% to 92%, with an average of 53.5% over the 44 games that were reported in *The Boston Globe*. The null hypothesis was that Bird’s shooting percentages resulted from 44 binomial distributions all having the same probability of success; that is, from 44 batches of independent coin flips with a “Larry Bird coin.”

It turned out that the null-hypothetical binomial model fit these data fairly well. That is, there was not enough variability in these data to reject the null hypothesis that each shot was exactly analogous to a flip of a coin. Here, however, it was not clear whether the data failed to reject the null

hypothesis because the null hypothesis was true or because the sample size was too small. To check this, an alternative to the binomial could be constructed and the posterior odds of the binomial calculated. A particularly simple alternative is to assume that the binomial parameters vary across games according to a beta distribution. In addition, Kass collected more data, using Vinnie Johnson, the player who at that time seemed to be the most commonly cited example of a “streak shooter,” as opposed to Larry Bird, who was viewed as a relatively stable shooter. Kass obtained Johnson’s shooting data for 380 games during the 1985–1989 seasons.

Bayes factors allow easy comparison of nonnested models and of irregular models. This is illustrated in the next application.

Application 3: Ozone Exceedances

As described by Smith (1989), ground-level ozone is a topic of concern because high levels of ozone indicate that the air is polluted. U.S. standards specify that a threshold level be exceeded on no more than three days in any three-year period. In a number of U.S. cities, including Houston, Texas, this standard is far from being met, and the task of regulatory bodies such as the Texas Air Control Board is to introduce measures to reduce the frequency and level of high exceedances.

The question of whether ozone exceedances were decreasing in Houston may be examined using the times of occurrence of clusters of exceedances above the threshold of 16 parts per 100 million (pphm) between 1974 and 1986, de-seasonalized; the data were given by Raftery (1989, Table 1). Exceedances may be modeled as a Poisson process with a rate that is perhaps varying over time. Raftery (1989) considered three models for the rate that represent (a) no change, (b) a gradual decrease, and (c) an abrupt decrease. There are several interesting features of this application. Models (b) and (c) are not nested. In addition, model (c) is not a regular statistical model in that it has a highly discontinuous likelihood function. As a result, frequentist methods for testing (a) against (c) are complex to develop, whereas Bayes factors are fairly simple to calculate. We know of no frequentist way of testing (b) against (c). The analysis suggested a more precise mechanism and, potentially, a more developed statistical model, illustrating that Bayes factors are not restricted to the comparison of previously formulated hypotheses but are also useful for guiding an evolving model-building process.

In the foregoing applications, simple plausible models were used to represent competing hypotheses. In each case, though, one could phrase the question differently by asking which of the models fits the data better. That is, the problem of testing a hypothesis was identified with one of model selection. But the connotation of the phrase “model selection” differs according to whether the model comes from underlying theory or from looking at the data. If the latter, then Bayes factors can still be used, but choosing a model may no longer be the primary goal.

This arises, for example, when one is interested in the effect of one variable on another and where there are several other possible covariates to be included in a regression equation. One

often is not sure of having chosen exactly the right set of covariates, and this is a source of uncertainty that should be taken into account. Similarly, other functional or distributional assumptions may lead to different estimates of quantities of interest, and again, one would like to take account of uncertainty about the assumptions within the estimation process. The Bayesian approach allows this to be done in a natural way by averaging over the candidate models with their posterior probabilities as weights. The following application raises this issue, where the competing models correspond to different link functions in a generalized linear model.

Application 4: Educational Transitions

Do social class background, ability, and type of school attended affect educational attainment? These questions were addressed in the context of Ireland by Greaney and Kelleghan (1984), who concluded that the Irish educational system is approaching the meritocratic ideal. By this they meant that progress within the system is determined largely by educationally relevant attributes and not by other, educationally irrelevant, attributes such as social class. One question of particular policy interest is whether and to what extent students in vocational (nonacademic second-level) schools drop out of school earlier than their counterparts in secondary (academic second-level) schools with the same ability and social class origin.

The questions may be addressed by reanalyzing the longitudinal data of Greaney and Kelleghan (1984) using models based on logistic regression. They then reduce to questions about the presence or absence of the effects of interest in the regression model and their size. Questions about the size of these effects have to be answered in the presence of several competing models (in this case different link functions) between which the data do not distinguish clearly but which nevertheless yield very different results. It is essential to take account of this model uncertainty.

In the final application, simple approximate methods were applied in a modestly complicated setting (i.e., repeated-measures logistic regression with alternative predictors). It was desired to check the modeling assumptions and the accuracy of the asymptotics used. This raised problems of computation and the determination of prior distributions.

Application 5: Human Working Memory Failure in Computer-Based Tasks

In human–computer interaction it is of interest to characterize tasks that tend to lead to human error, so that procedures may be written to avoid them. Carlin, Kass, Lerch, and Huguenard (1992) have described one such effort in an experimental study involving the data base management system SQL, which presented subjects with query tasks of varying complexity (such as “find all customers having an outstanding invoice of more than \$200”). Understanding of errors is guided by cognitive psychological theory concerning overload of what is now called “working memory” (a concept that has replaced and refined what used to be called short-term memory). Thus the notion is that errors tend to occur when human working memory is overloaded because of ex-

cessive demands in the task. At issue in this study was whether overload tended to be due to (a) the number of conditions in the query or (b) the complexity of the query. To determine which characterization led to better predictions of error rates, two alternative statistical models were constructed based on alternative predictor variables, either the number of conditions or a pair of measures of query complexity.

There were 20 experimental subjects, each of whom was given 50 tasks to complete. Half the subjects received a "cue," which was supposed to remind them of the essential feature of the problem and to prevent them from making the error under study. Carlin et al. (1992) used a logistic regression model with random subject ability effects (a single number for each subject) nested within the cue effect, together with the respective predictor variables. The two competing models had four parameters in common; the number-of-conditions model had one additional parameter (the coefficient of the predictor), whereas the query-complexity model had two additional parameters (the coefficients of the two predictors). Evaluation of the likelihood in either model required the numerical calculation of an integral due to the random subject effects.

3. BAYES FACTORS

3.1 Definition

We begin with data \mathbf{D} , assumed to have arisen under one of the two hypotheses H_1 and H_2 according to a probability density $\text{pr}(\mathbf{D}|H_1)$ or $\text{pr}(\mathbf{D}|H_2)$. Given a priori probabilities $\text{pr}(H_1)$ and $\text{pr}(H_2) = 1 - \text{pr}(H_1)$, the data produce a posteriori probabilities $\text{pr}(H_1|\mathbf{D})$ and $\text{pr}(H_2|\mathbf{D}) = 1 - \text{pr}(H_1|\mathbf{D})$. Because any prior opinion gets transformed to a posterior opinion through consideration of the data, the transformation itself represents the evidence provided by the data. In fact, the same transformation is used to obtain the posterior probability, regardless of the prior probability. Once we convert to the odds scale (odds = probability/(1 - probability)), the transformation takes a simple form. From Bayes's theorem, we obtain

$$\text{pr}(H_k|\mathbf{D}) = \frac{\text{pr}(\mathbf{D}|H_k)\text{pr}(H_k)}{\text{pr}(\mathbf{D}|H_1)\text{pr}(H_1) + \text{pr}(\mathbf{D}|H_2)\text{pr}(H_2)} \quad (k = 1, 2),$$

so that

$$\frac{\text{pr}(H_1|\mathbf{D})}{\text{pr}(H_2|\mathbf{D})} = \frac{\text{pr}(\mathbf{D}|H_1)\text{pr}(H_1)}{\text{pr}(\mathbf{D}|H_2)\text{pr}(H_2)},$$

and the transformation is simply multiplication by

$$B_{12} = \frac{\text{pr}(\mathbf{D}|H_1)}{\text{pr}(\mathbf{D}|H_2)}, \quad (1)$$

which is the *Bayes factor*. Thus, in words,

$$\text{posterior odds} = \text{Bayes factor} \times \text{prior odds},$$

and the Bayes factor is the ratio of the posterior odds of H_1 to its prior odds, regardless of the value of the prior odds. (The terminology is apparently due to Good 1958, who at-

tributed the method to Turing in addition to, and independently of, Jeffreys at about the same time; see Good 1983.) When the hypotheses H_1 and H_2 are equally probable a priori so that $\text{pr}(H_1) = \text{pr}(H_2) = .5$, the Bayes factor is equal to the posterior odds in favor of H_1 . The two hypotheses may well not be equally likely a priori, however.

In the simplest case, when the two hypotheses are single distributions with no free parameters (the case of "simple versus simple" testing), B_{12} is the likelihood ratio. In other cases, when there are unknown parameters under either or both of the hypotheses, the Bayes factor is still given by (1), and, in a sense, it continues to have the form of a likelihood ratio. Then, however, the densities $\text{pr}(\mathbf{D}|H_k)$ ($k = 1, 2$) are obtained by *integrating* (not maximizing) over the parameter space, so that in Equation (1),

$$\text{pr}(\mathbf{D}|H_k) = \int \text{pr}(\mathbf{D}|\theta_k, H_k)\pi(\theta_k|H_k) d\theta_k, \quad (2)$$

where θ_k is the parameter under H_k , $\pi(\theta_k|H_k)$ is its prior density, and $\text{pr}(\mathbf{D}|\theta_k, H_k)$ is the probability density of \mathbf{D} given the value of θ_k , or the likelihood function of θ . Here θ_k may be a vector, and in what follows we will denote its dimension by d_k .

The prior distributions $\pi(\theta_k|H_k)$ ($k = 1, 2$) are necessary. This may be considered both good and bad. Good, because it is a way of including other information about the values of the parameters (as in Application 1). Bad, because these prior densities may be hard to set when there is no such information. We discuss the problem of setting priors and assessing sensitivity to the choices in Section 5.

The quantity $\text{pr}(\mathbf{D}|H_k)$ given by Equation (2) is the marginal probability of the data, because it is obtained by integrating the joint density of (\mathbf{D}, θ_k) given \mathbf{D} over θ_k . It is also the predictive probability of the data; that is, the probability of seeing the data that actually were observed, calculated *before* any data became available. It is also sometimes called a marginal likelihood, or an integrated likelihood. Note that, as in computing the likelihood ratio statistic but unlike in some other applications of likelihood, all constants appearing in the definition of the likelihood $\text{pr}(\mathbf{D}|\theta_k, H_k)$ must be retained when computing B_{12} . In fact, B_{12} is closely related to the likelihood ratio statistic, in which the parameters θ_k are eliminated by maximization rather than by integration. We discuss this relationship in Sections 4.1.2 and 4.1.3.

Other notations are often used for the Bayes factor. When many hypotheses are involved, we will write B_{jk} as the Bayes factor for H_j against H_k . Often, one of the hypotheses is considered the null and thus is denoted by H_0 . In this case if there is only one alternative, it will be denoted by H_1 , so that, putting $\text{pr}(\mathbf{D}|H_0)$ in the numerator, the Bayes factor becomes B_{01} . In this situation Jeffreys instead used K . When comparing results with standard likelihood ratio tests, it is convenient to instead put the null hypothesis in the denominator of (1) and thus use B_{10} as the Bayes factor.

For the usual (non-Bayesian) large-sample distribution theory of likelihood ratio tests to be applicable, the null must be nested within the alternative. That is, there must be some parameterization under H_1 of the form $\theta = (\beta, \psi)$ such that H_0 is obtained from H_1 when $\psi = \psi_0$ for some ψ_0 . Here both

β and ψ may be vectors. Although the expression (1) does not require the models to be nested, the case of nested models is of special interest in the Bayesian approach as well, and we will refer to it frequently in what follows.

3.2 Interpretation

The Bayes factor is a summary of the evidence provided by the data in favor of one scientific theory, represented by a statistical model, as opposed to another. Jeffreys (1961, app. B) suggested interpreting B_{10} in half-units on the \log_{10} scale. Pooling two of his categories together for simplification, we have:

$\log_{10}(B_{10})$	B_{10}	Evidence against H_0
0 to 1/2	1 to 3.2	Not worth more than a bare mention
1/2 to 1	3.2 to 10	Substantial
1 to 2	10 to 100	Strong
>2	>100	Decisive

Probability itself provides a meaningful scale defined by betting, and so these categories are not a calibration of the Bayes factor, but rather a rough descriptive statement about standards of evidence in scientific investigation. We speak here in terms of B_{10} , because weighing evidence *against* a null hypothesis is more familiar, but Bayes factors can equally well provide evidence *in favor* of a null hypothesis. Of course, the interpretation may depend on the context. For example, Evett (1991) has argued that for forensic evidence alone to be conclusive in a criminal trial, one would require posterior odds for H_1 (guilt) against H_0 (innocence) of at least 1,000 rather than the 100 suggested by Jeffreys.

It can be useful to consider twice the natural logarithm of the Bayes factor, which is on the same scale as the familiar deviance and likelihood ratio test statistics. Rounding and using 20 rather than 10 as the requirement for strong evidence, we then obtain a slight modification:

$2 \log_e(B_{10})$	(B_{10})	Evidence against H_0
0 to 2	1 to 3	Not worth more than a bare mention
2 to 6	3 to 20	Positive
6 to 10	20 to 150	Strong
>10	>150	Very strong

From our own experience, these categories seem to furnish appropriate guidelines.

The logarithm of the marginal probability of the data may also be viewed as a predictive score. This is of interest, because it leads to an interpretation of the Bayes factor that does not depend on viewing one of the models as “true.” Suppose that $\mathbf{D} = \{y_1, \dots, y_n\}$ and that for each i , we form a predictive distribution $\widehat{\text{pr}}_i(\cdot)$ of y_i given the already available data $\{y_1, \dots, y_{i-1}\}$. We use the logarithmic scoring rule, $\log \widehat{\text{pr}}_i(y_i)$ (Good 1952), to assess performance. Then the overall score of any rule that generates such predictive distributions is $\text{LS} = \sum_i \log \widehat{\text{pr}}_i(y_i)$. In particular, if the prediction rule is derived from the model H_k (i.e., likelihood and prior), then $\log \text{pr}(\mathbf{D}|H_k) = \sum_i \log \text{pr}(y_i | y_{i-1}, \dots, y_1,$

$H_k) = \text{LS}_k$. It follows that the log Bayes factor is $\log B_{10} = \text{LS}_1 - \text{LS}_0$; that is, the difference in predictive scores. Thus the Bayes factor can be viewed as measuring the relative success of H_1 and H_0 at predicting the data. This is related to *prequential analysis* (Dawid 1984) and also to *stochastic complexity* (Rissanen 1987); the connections were discussed by Dawid (1992) and Hartigan (1992). Good (1985), and in many other publications, has referred to the log Bayes factor as the “weight of evidence.”

4. CALCULATING BAYES FACTORS

In some elementary cases the integral (2), which we will rewrite in this section as

$$I = \int \text{pr}(\mathbf{D}|\theta, H)\pi(\theta|H) d\theta, \tag{3}$$

may be evaluated analytically. More often, it is intractable and thus must be computed by numerical methods. But most available software developed by numerical analysts is generally so inefficient for these integrals that it is of little use. One reason is that when sample sizes are moderate or large, the integrand becomes highly peaked around its maximum, which may be found by other techniques, and quadrature methods that do not begin with knowledge of the maximum are likely to have difficulty finding the region where the integrand mass is accumulating. A second reason is that some problems are of high dimension. In this case Monte Carlo methods may be used, but these too need to be adapted to the statistical context. A review of various numerical integration strategies for evaluating the integral in (3) is provided by Evans and Swartz (1995).

Exact analytic evaluation of the integral (3) is possible for exponential family distributions with conjugate priors, including normal linear models (DeGroot 1970, chap. 9; Zellner 1971, chap. 10).

4.1 Asymptotic Approximation

4.1.1 Laplace’s Method. A useful approximation to the marginal density of the data as given by (3) is obtained by assuming that the posterior density, which is proportional to $(\text{pr}(\mathbf{D}|\theta, H)\pi(\theta|H))$, is highly peaked about its maximum $\tilde{\theta}$, which is the posterior mode. This will usually be the case if the likelihood function $\text{pr}(\mathbf{D}|\theta, H)$ is highly peaked near its maximum $\tilde{\theta}$, which will be the case for large samples. Let $\tilde{l}(\theta) = \log(\text{pr}(\mathbf{D}|\theta, H)\pi(\theta|H))$. Expanding $\tilde{l}(\theta)$ as a quadratic about $\tilde{\theta}$ and then exponentiating yields an approximation to $(\text{pr}(\mathbf{D}|\theta, H)\pi(\theta|H))$ that has the form of a normal density with mean $\tilde{\theta}$ and covariance matrix $\tilde{\Sigma} = (-\mathbf{D}^2\tilde{l}(\tilde{\theta}))^{-1}$, where $\mathbf{D}^2\tilde{l}(\tilde{\theta})$ is the Hessian matrix of second derivatives. Integrating this approximation yields

$$\hat{I} = (2\pi)^{d/2} |\tilde{\Sigma}|^{1/2} \text{pr}(\mathbf{D}|\tilde{\theta}, H)\pi(\tilde{\theta}|H), \tag{4}$$

where d is the dimension of θ .

This is Laplace’s method of approximation (de Bruijn 1970, sec. 4.4; Tierney and Kadane 1986). For many problems in which the sample size n is moderate, it produces answers well within the accuracy required for drawing conclusions according to the scheme of Section 3.2. Formally,

under conditions spelled out by Kass, Tierney, and Kadane (1990), as $n \rightarrow \infty$, $I = \hat{I}(1 + O(n^{-1}))$; that is, the relative error is $O(n^{-1})$. Thus when Laplace's method is applied to both the numerator and denominator of B_{10} in (1), the resulting approximation also has relative error of order $O(n^{-1})$.

The accuracy of Laplace's method in this context has been examined by Kass and Vaidyanathan (1992), Rosenkranz (1992), and Raftery (1993c) and is mentioned in Applications 4 and 5. In general, the method provides adequate approximations for well-behaved problems (those in which the likelihood functions are not grossly nonnormal) of modest dimensionality. We are unable to be much more specific than this, though our rough feeling is that samples of size less than $5d$ are worrisome (with d being the dimension of θ), whereas those of size greater than $20d$ are large enough for the method to work well in most familiar problems, provided that a reasonably good parameterization is used. (See Slate 1994 for a detailed discussion of sample sizes required to obtain posterior normality, which would guarantee accuracy of Laplace's method, for various parameterizations of exponential families.) Software for applying Laplace's method is available in *S* and in *LispStat* (Raftery 1993c; Tierney 1989, 1990). Approximations of this kind have been used by many authors, notably Jeffreys (1961), Lindley (1961), Mosteller and Wallace (1964), and Leonard (1982).

4.1.2 Variants on Laplace's Method. Laplace's method may be applied in alternative forms by omitting part of the integrand from the exponent when performing the expansion. (For the general formulation, see Kass and Vaidyanathan 1992, which followed Tierney, Kass, and Kadane 1989 and Mosteller and Wallace 1964, Sec. 4.6). An important variant on (4) is

$$\hat{I}_{MLE} = (2\pi)^{d/2} |\hat{\Sigma}|^{1/2} \text{pr}(\mathbf{D}|\hat{\theta}, H) \pi(\hat{\theta}|H) \quad (5)$$

where $\hat{\Sigma}^{-1}$ is the observed information matrix; that is, the negative Hessian matrix of the log-likelihood evaluated at the maximum likelihood estimator (MLE) $\hat{\theta}$. This approximation again has relative error of order $O(n^{-1})$. Although it is likely to be less accurate than (4) when the prior is somewhat informative relative to the likelihood, it has the advantage that it is easily computed from any statistical software package that reports the MLE, the observed information matrix (or its inverse), and the value of the maximized likelihood.

Some software packages calculate the expected information matrix (i.e., the usual Fisher information matrix), or its inverse, as the asymptotic covariance matrix rather than the observed information matrix. The inverse of the expected information matrix may be used in place of $\hat{\Sigma}$ in (5). The resulting approximation has a larger asymptotic relative error, of order $O(n^{-1/2})$, but it remains sufficiently accurate to be of use in many problems.

Now suppose that we have nested hypotheses, with parameter (β, ψ) having prior $\pi(\beta, \psi|H_1)$ under H_1 and then $H_0: \psi = \psi_0$ with prior $\pi(\beta|H_0)$. When (5) is applied, we obtain

$$2 \log B_{10} \approx \Lambda + \log |\hat{\Sigma}_1| - \log |\hat{\Sigma}_0| + \log \pi(\hat{\beta}, \hat{\psi}|H_1) - \log \pi(\hat{\beta}^*, |H_0) + (d_1 - d_0) \log(2\pi), \quad (6)$$

where $\Lambda = 2(\log \text{pr}(\mathbf{D} | (\hat{\beta}, \hat{\psi}), H_1) - \log \text{pr}(\mathbf{D} | \hat{\beta}, H_0))$ is

the log-likelihood ratio statistic having degrees of freedom $(d_1 - d_0)$ and $\hat{\beta}^*$ denotes the MLE under H_0 . Again, either observed or expected information may be used in computing the covariance matrices $\hat{\Sigma}_k$; then (6) has relative error of order $O(n^{-1})$ or $O(n^{-1/2})$. Jeffreys (1961, sec. 5.31) gave an approximation for the case where $(d_1 - d_0) = 1$ in essentially the form of (6) but with expected rather than observed information, and Chow (1981) extended Jeffreys' result to higher dimensions. Hsiao (1994) investigated alternative expansions for the case in which $\hat{\psi} = \psi_0$ lies at the boundary of the parameter space.

Raftery (1993c) suggested approximating $\tilde{\theta}_k$ by a single Newton step starting from $\hat{\theta}_k$ and substituting the result into Equation (4), which yields the approximation

$$2 \log B_{10} \approx \Lambda + (E_1 - E_0). \quad (7)$$

In Equation (7),

$$E_k = 2\lambda_k(\hat{\theta}_k) + \lambda'_k(\hat{\theta}_k)^T (F_k + G_k)^{-1} \{2 - F_k(F_k + G_k)^{-1}\} \times \lambda'_k(\hat{\theta}_k) - \log |F_k + G_k| + d_k \log(2\pi), \quad (8)$$

where $\text{var}[\theta_k|H_k] = W_k$, $F_k = \hat{\Sigma}_k^{-1}$, $G_k = W_k^{-1}$, $\lambda_k(\theta_k) = \log \text{pr}(\theta_k|H_k)$ is the log-prior density, and $\lambda'_k(\hat{\theta}_k)$ is the d_k vector of derivatives of $\lambda_k(\theta_k)$ with respect to the elements of θ_k ($k = 0, 1$). This often improves on the approximation of (5), but it does not require any additional information. It is implemented for generalized linear models in the *GLIB* software; see Section 9.

4.1.3 The Schwarz Criterion. It is possible to avoid the introduction of the prior densities $\pi_k(\theta_k|H_k)$ in (1) by using $S = \log \text{pr}(\mathbf{D}|\hat{\theta}_1, H_1) - \log \text{pr}(\mathbf{D}|\hat{\theta}_2, H_2)$

$$- \frac{1}{2} (d_1 - d_2) \log(n),$$

where $\hat{\theta}_k$ is the MLE under H_k , d_k is the dimension of θ_k , and n is the sample size. As $n \rightarrow \infty$, this quantity, often called the Schwarz criterion, satisfies

$$\frac{S - \log B_{12}}{\log B_{12}} \rightarrow 0 \quad (9)$$

and thus may be viewed as a rough approximation to the logarithm of the Bayes factor. Minus twice the Schwarz criterion is often called the Bayesian information criterion (BIC); sometimes an arbitrary constant is added. See Section 8.3 for additional discussion.

In contrast to the approximations furnished by (4) or (5), the relative error of $\exp(S)$ in approximating B_{12} is generally $O(1)$. Thus even for very large samples, it does not produce the correct value. On the other hand, keeping in mind the rough interpretation of B_{12} on the logarithmic scale suggested in Section 3.2, Equation (9) shows that in large samples the Schwarz criterion should provide a reasonable indication of the evidence.

The Schwarz criterion is appealing in that it can be applied as a standard procedure even when the priors $\pi(\theta_k|H_k)$ are hard to set precisely. In this sense it provides an often-useful reference procedure for scientific reporting. Also in its favor is the following. For independent and identically distributed

sampling in nested models with $H_0: \psi = \psi_0$, it is reasonable to take the prior on ψ to be $N(\psi_0, W)$ with $|W| = |\mathbf{J}(\psi_0)|^{-1}$, where $\mathbf{J}(\psi)$ is the Fisher information in the parameter ψ . This says that the amount of information in the prior is equal to the amount of information in one observation. In this case, under certain conditions (i.e., the assumptions of null orthogonality and local alternative mentioned in Sec. 5.2), the Schwarz criterion furnishes an order $O(n^{-1/2})$ approximation to $\log B_{12}$ (Kass and Wasserman 1995). Thus if one were willing to use this prior as a “reference prior” suitable for standardized reporting, then the Schwarz criterion would be a reasonably good approximation to the log of the Bayes factor. As Kass and Wasserman (1995) noted, Jeffreys’s method was essentially to use $S + c$ as an approximation to the log of the Bayes factor, with c being a constant determined by substituting a Cauchy prior in place of the normal prior on ψ .

The sample size n in the definition of S needs to be figured carefully. It is apparent from the derivation of S (e.g., Kass 1993), that n should be the rate at which the Hessian matrix of the log-likelihood function grows; thus n becomes the number of data values contributing to the summation that appears in the formula for the Hessian. For instance, as Raftery (1986a) pointed out, in the case of log-linear models for contingency tables, n is the sum of the counts, not of the number of cells. Similarly, in models for binomial responses, n is the sum of the denominators, not the number of responses or the number of successes. In survival analysis, Raftery, Madigan, and Volinsky (1995) have taken n to be the number of *uncensored* observations; that is, of deaths. In two-stage hierarchical models the situation can be subtle and the appropriate form of S depends on which parameters are being tested.

Schwarz (1978) gave a rigorous derivation of the criterion for linear subfamilies of exponential families and Houghton (1988) extended Schwarz’s result to curved exponential families. Heuristic derivations of (9) are quite easy: one applies (5) and then neglects constant-order terms. It is apparent from arguments such as those of Kass et al. (1990) that (9) holds much more generally than in the restricted setting of curved exponential families (e.g., Kass 1993, Katz 1981, Leonard 1982, Raftery 1986a). But rigorously demonstrating the assumptions for the validity of Laplace’s method seems to be enough of a chore that no general precisely formulated result has been published.

4.2 Simple Monte Carlo, Importance Sampling, and Gaussian Quadrature

Dropping the notational dependence on H_k , Equation (2) becomes

$$\text{pr}(\mathbf{D}) = \int \text{pr}(\mathbf{D}|\theta)\pi(\theta) d\theta.$$

The simplest Monte Carlo integration estimate of this is

$$\widehat{\text{pr}}_1(\mathbf{D}) = \frac{1}{m} \sum_{i=1}^m \text{pr}(\mathbf{D}|\theta^{(i)}),$$

where $\{\theta^{(i)}: i = 1, \dots, m\}$ is a sample from the prior distribution; this is the average of the likelihoods of the sampled parameter values (e.g., Hammersley and Handscomb 1964).

This possibility was mentioned by Raftery and Banfield (1991) and was investigated in some detail in particular cases by McCulloch and Rossi (1991). A major difficulty with $\widehat{\text{pr}}_1(\mathbf{D})$ is that most of the $\theta^{(i)}$ have small likelihood values if the posterior is concentrated relative to the prior, so that the simulation process will be quite inefficient. Thus the estimate is dominated by a few large values of the likelihood, and so the variance of $\widehat{\text{pr}}_1(\mathbf{D})$ is large and its convergence to a Gaussian distribution is slow. These problems were apparent in the examples studied in detail by McCulloch and Rossi (1991).

The precision of simple Monte Carlo integration can be improved by importance sampling. This consists of generating a sample $\{\theta^{(i)}: i = 1, \dots, m\}$ from a density $\pi^*(\theta)$. Under quite general conditions, a simulation-consistent estimate of I is

$$\hat{I} = \frac{\sum_{i=1}^m w_i \text{pr}(\mathbf{D}|\theta^{(i)})}{\sum_{i=1}^m w_i}, \quad (10)$$

where $w_i = \pi(\theta^{(i)})/\pi^*(\theta^{(i)})$; the function $\pi^*(\theta)$ is known as the *importance sampling function*. (For general discussion and references, see Geweke 1989.)

A more efficient scheme is based on adaptive Gaussian quadrature. Using well-established methods from the numerical analysis literature, Genz and Kass (1993) showed how integrals that are peaked around a dominant mode may be evaluated. This approach is effective in such problems when the dimensionality of the parameter space is modest (roughly, less than 9).

4.3 Simulating from the Posterior

Several methods are now available for simulating from posterior distributions. In the simplest cases these include direct simulation and rejection sampling. In more complex cases, Markov chain Monte Carlo (MCMC) methods, particularly the Metropolis–Hastings algorithm and the Gibbs sampler, provide a general recipe (e.g., Smith and Roberts 1993 and references therein). Another fairly general recipe is the weighted likelihood bootstrap (Newton and Raftery 1994).

Any of these methods gives us a sample approximately drawn from the posterior density $\pi^*(\theta) = \text{pr}(\theta|\mathbf{D}) = \text{pr}(\mathbf{D}|\theta)\pi(\theta)/\text{pr}(\mathbf{D})$. Substituting into Equation (10) yields as an estimate for $\text{pr}(\mathbf{D})$,

$$\widehat{\text{pr}}_2(\mathbf{D}) = \left\{ \frac{1}{m} \sum_{i=1}^m \text{pr}(\mathbf{D}|\theta^{(i)})^{-1} \right\}^{-1}, \quad (11)$$

the harmonic mean of the likelihood values (Newton and Raftery 1994). This converges almost surely to the correct value, $\text{pr}(\mathbf{D})$, as $m \rightarrow \infty$, but it does not generally satisfy a Gaussian central limit theorem. This manifests itself by the occasional occurrence of a value of $\theta^{(i)}$ with small likelihood and hence large effect on the final result. But it is very easy to calculate, and experience to date suggests that although it is indeed unstable, it often gives results that are accurate

enough for interpretation on the logarithmic scale of Section 3.2 (Carlin and Chib 1993; Raftery 1994; Rosenkranz 1992).

Several modifications of the harmonic mean estimator (11) have been suggested to get around its instability. Newton and Raftery (1994) suggested using as importance sampling function in (10) a mixture of the prior and posterior densities, $\pi^*(\theta) = \delta\pi(\theta) + (1 - \delta)\text{pr}(\theta|\mathbf{D})$, where $0 < \delta < 1$. The resulting estimator, $\widehat{\text{pr}}_3(\mathbf{D})$, has the efficiency of $\widehat{\text{pr}}_2(\mathbf{D})$ because it is based on many values of θ with high likelihood, but avoids its instability and does satisfy a Gaussian central limit theorem. But it has the irksome aspect that one must simulate from the prior as well as the posterior. This may be avoided by simulating all m values from the posterior distribution and imagining that a further $\delta m/(1 - \delta)$ values of θ are drawn from the prior, all of them with likelihoods $\text{pr}(\mathbf{D}|\theta^{(i)})$ equal to their expected value $\text{pr}(\mathbf{D})$. The resulting estimator, $\widehat{\text{pr}}_4(\mathbf{D})$, may be evaluated using a simple iterative scheme.

A simple modification of (11) is to instead calculate

$$\widehat{\text{pr}}_5(\mathbf{D}) = \frac{1}{m} \sum_{i=1}^m f(\theta^{(i)}) / (\text{pr}(\mathbf{D}|\theta^{(i)})\pi(\theta^{(i)})), \quad (12)$$

where $f(\cdot)$ is any d -dimensional probability density. This was mentioned by Gelfand and Dey (1994). It is an unbiased and consistent estimator of the marginal likelihood $\text{pr}(\mathbf{D})$, and satisfies a Gaussian central limit theorem if the tails of $f(\cdot)$ are thin enough, specifically if $\int \{f(\theta)^2 / (\text{pr}(\mathbf{D}|\theta)\pi(\theta))\} d\theta < \infty$. High efficiency would seem most likely to result if $f(\theta)$ were roughly proportional to $\text{pr}(\mathbf{D}|\theta)$. The very limited experience to date indicates that for low-dimensional problems with a good choice of f , $\widehat{\text{pr}}_5(\mathbf{D})$ can be very accurate. For high-dimensional problems, however, it may be hard to find an appropriate f , and the results can be very poor. Meng and Wong (1993) considered an optimal choice of f and showed how this could be computed iteratively from an initial guess. [Their framework actually applies to the computation of a ratio of integrals of the form (3).] Their approach appears promising but has not yet been extensively tested.

Raftery (1995a) suggested what he called the ‘‘Laplace-Metropolis’’ estimator of $\text{pr}(\mathbf{D})$, obtained by using the posterior simulation output to estimate the quantities needed to compute the Laplace approximation (4), namely the posterior mode, $\tilde{\theta}$, and minus the inverse Hessian at the posterior mode, $\tilde{\Sigma}$. The posterior mode can be estimated as the $\theta^{(i)}$ that maximizes $(\text{pr}(\mathbf{D}|\theta^{(i)})\pi(\theta^{(i)}))$. This requires computing the likelihood for each simulated $\theta^{(i)}$; if this takes too much computer time, then an alternative is to use the multivariate or componentwise posterior median or to estimate the posterior mode by nonparametric density estimation. The matrix $\tilde{\Sigma}$ can be estimated by the estimated posterior covariance matrix; it is wise to use a high-breakdown point robust estimator. The resulting estimator has performed well in numerical experiments (Lewis and Raftery 1994). A similar combination of simulation with Laplace’s method was suggested by Kass and Wasserman (1992a), who provided a correction term as well. From recent unpublished studies we have conducted with T. J. DiCiccio and L. Wasserman, the latter method appears to be quite promising.

In the case of nested models in which the priors satisfy (13) of Section 5.1 and β and ψ are independent, then

$$B_{10} = \frac{\pi(\psi_0|H_1)}{\text{pr}(\psi_0|\mathbf{D}, H_1)}.$$

(Actually, instead of a priori independence of β and ψ , it is enough to have $\pi(\beta|H_0) = \pi(\beta|\psi = \psi_0, H_1)$.) This result is apparently due to L. J. Savage and was called the ‘‘Savage density ratio method’’ by Dickey (1971). It has been exploited in several applications (e.g., McCulloch and Rossi 1991). Verdinelli and Wasserman (1993a) pointed out that a method of generating marginal posterior densities discussed by Chen (1992) may then be applied advantageously in this case. This leads to a method similar to $\widehat{\text{pr}}_5(\mathbf{D})$, and the choice of f is again crucial and difficult.

Carlin and Chib (1993) suggested including a model indicator variable in the MCMC scheme and defining ‘‘pseudo-priors’’ for $(\theta_1|H_2)$ and $(\theta_2|H_1)$. This involves designing and running a special MCMC algorithm to calculate Bayes factors. Similar suggestions have been made by Carlin and Polson (1991) and George and McCulloch (1993).

An alternative approach is available when many parameters, z , are present in all the models considered. These might be missing data, the values of a random effect in a hierarchical model, or ‘‘latent data’’ chosen in such a way that the ‘‘complete data likelihood’’ $\text{pr}(\mathbf{D}, z|\theta)$ has a simple form (Tanner and Wong 1987). Then the Bayes factor can be simulation-consistently estimated by the average of the quantities $B_{10}(z^{(i)}) = \text{pr}(\mathbf{D}, z^{(i)}|H_1) / \text{pr}(\mathbf{D}, z^{(i)}|H_0)$, where the $z^{(i)}$ are simulated from the posterior distribution of z under H_0 ; this is often possible using MCMC methods. The $B_{10}(z^{(i)})$ are then often easy to calculate or at least to approximate fairly well; for example, using the Laplace method. When z is present in H_1 but not in H_0 , we again recover the harmonic mean estimator of $\text{pr}(\mathbf{D}|H_1)$ (Raftery 1993a). This is related to previous work of Thompson and Wijsman (1990) on the calculation of likelihood ratios.

Finally, some general methods of calculating I have been considered in the statistical physics literature under the name ‘‘free energy estimation.’’ The approaches are not automatic and require analytical effort to tailor them to statistical applications; however, they may be of use in certain problems (See Neal 1992 for references.)

4.4 Comparison of Methods

The different methods for calculating and approximating Bayes factors have been compared by Rosenkrantz (1992) in the contexts of normal models and of Poisson-gamma models for counts with independent variables, unobserved heterogeneity, and outliers, and by Raftery (1993c) in the context of generalized linear models.

Exact analytic evaluation is best because it is most accurate and usually also most efficient computationally, but it is feasible only for a narrow class of models. The Laplace method for integrals yields accurate approximations and is usually quite computationally efficient. In particular, the approximation using the posterior mode in (4) or its one-step approximation given by Equations (7) and (8) can be very accurate. The latter is easy to compute using the output from

standard software. For cases of modest dimensionality, the adaptive quadrature method of Genz and Kass (1993) is effective. The Monte Carlo integration and importance sampling methods are less precise and more computationally demanding, but they may well be the only ones available in complex models. The methods using MCMC seem promising but have not yet been applied in many demanding problems. In addition, they can require large numbers of likelihood function evaluations, which in some cases is itself difficult.

The Schwarz criterion is the easiest approximation to compute and has the advantage of not requiring that the user specify the prior distributions. It is suitable for communicating results even when more detailed calculations are made using other methods (as in Application 5). As an asymptotic approximation it is admittedly rough, but, as long as the number of degrees of freedom involved in the comparison is reasonably small relative to sample size, it does not seem to be grossly misleading in a qualitative sense. It can be very poor, however, when the number of degrees of freedom involved in the comparison is large and the prior is very different from that for which the approximation is best; McCulloch and Rossi (1991) gave an example of this with 115 degrees of freedom.

5. THE CHOICE OF PRIORS

In order to compute a Bayes factor, the prior distributions $\pi(\theta_k|H_k)$ on the parameters of each model must be specified. Sometimes, as in Application 1, there are closely related data with which to construct priors. This nicest situation is rare, however. More often, as in Application 5, some combination of relevant data, information from the literature, and rough guesses must be used. In that case there will be doubts about the accuracy of the prior distribution. Thus a first concern is how to choose prior distributions to represent the available information, but once this is done, an important issue is the sensitivity of the Bayes factor to the choices of priors.

The easiest way to deal with the problem of prior choice is to ignore it and simply use the Schwarz criterion or Jeffreys's variant of it (see Sec. 4.1.3). Although this will lead to appropriate conclusions in sufficiently large samples, there is not much available guidance as to the operational meaning of "sufficiently large." Also, in contrast with Bayesian point estimates such as the posterior mean, the Bayes factor does tend to be sensitive to the choices of priors on the model parameters. We discuss informative prior selection and the problem of sensitivity in Section 5.1 and review results and methods for assessing sensitivity over a range of priors in Section 5.2. In Section 5.3 we consider the use of improper priors and note the resulting difficulties.

Sensitivity analysis concerns distributional forms for models $\text{pr}(\mathbf{D}|\theta_k, H_k)$ as well as priors. When alternatives are introduced (e.g., the Student's t distribution in place of the normal), Bayes factors may be used to determine which best fits the data. One may also assess the influence of individual data values by computing the Bayes factor after omitting each observation in turn (Pettit and Young 1990). Asymptotic approximation makes the "leave-one-out" diagnostic approach easy; see Application 5.

5.1 Prior Information

The problem of determining a prior distribution from available information appears throughout Bayesian inference. The information may come from other data or from the subjective knowledge of experts, or it may be too hard to express as a prior distribution. Even when there are other data, judgment must be used because a distribution must be chosen (as in Application 1), but this is familiar. The use of subjective opinion is different. For the most part, in both of these cases the formulation of priors is problem-specific and there is not much general methodology for it. There is, however, a modest literature on "eliciting" probabilistic information from individuals. One psychologically appealing device is the method of imaginary observations (Good 1950). Kadane, Dickey, Winkler, Smith, and Peters (1980) discussed a formal elicitation procedure in the context of linear regression; Garthwaite (1992) included more recent references. Formal representation of opinion may be useful for private analysis and in verifying that a simpler public analysis leads to appropriate conclusions, as in Application 5. Prior distributions sometimes become communally accepted after much research and discussion; this happened in the estimation of bowhead whale population size (Raftery and Zeh 1993).

In choosing priors, just as in choosing models for data distributions, simplifications are often made. This occurs notably when there are nested models in which a hypothesis $H_0: \psi = \psi_0$ is being tested in the presence of an additional parameter β . For example, it is often assumed that

$$\pi(\beta|H_0) = \int \pi(\beta, \psi|H_1) d\psi. \quad (13)$$

If in addition it is assumed that β and ψ are independent a priori under H_1 , then one needs only to choose one prior for β and another for ψ .

Sometimes instead of (13), it is assumed that $\pi(\beta|H_0) = \pi(\beta|\psi = \psi_0, H_1)$. But this prescription depends on the choice of parameterization used. If an alternative parameterization (ξ, λ) is used under H_1 , with H_0 then being specified by $\lambda = \lambda_0$, and if $\pi(\xi, \lambda|H_1)$ is obtained from $\pi(\beta, \psi|H_1)$ by the change-of-variables formula (introducing the appropriate Jacobian determinant), then it may happen that $\pi(\xi|H_0)$ is not obtained by a change of variables from $\pi(\beta|H_0)$. (In a personal communication, J. Dickey has noted that L. J. Savage discussed this as an instance of the "Borel paradox" in a 1963 lecture; see also Dickey 1985.) Variations on the marginal and conditional density methods have been discussed by McCulloch and Rossi (1993) and Verdinelli and Wasserman (1993a).

Again as in data modeling, simplifications involving priors should be considered carefully, because they may affect the results and yet may not be justified. Here testing is different from estimation. In frequentist theory, estimation and testing are complementary, but in the Bayesian approach, the problems are completely different. In testing, Bayesians put positive prior probability on models that represent hypotheses, whereas in estimation, continuous priors are used that would assign zero probability to special values such as $\psi = \psi_0$. In

estimation, priors are often picked for convenience, knowing that if the sample is fairly large, then the effect of the prior is small. In testing this is not so. As expression (5) shows, to order $O(n^{-1})$ the prior density may not be eliminated as the sample size increases, which contrasts with the case of estimation, where the MLE approximates the posterior mean to this order. As a result, the Bayes factor tends to be more sensitive to the choice of prior than the posterior probability of an interval (e.g., Kass 1993; Kass and Greenhouse 1989).

In fact, it may happen that conclusions based on estimation seem to contradict those from a Bayes factor (as illustrated in Kass and Greenhouse 1989). This occurs when an estimate, say of ψ , is found to be distant from a null-hypothetical value ψ_0 . In this case the data seem unlikely under H_0 , but if the Bayes factor turns out to be *in favor of* H_0 , then the data are *even more unlikely* under H_1 than they would have been under H_0 . (This is one cause of discrepancies between conclusions reached with Bayes factors and those reached using P values; see Sec. 8.2.) A consequence is that using a prior with a very large spread on ψ under H_1 in an effort to make it “noninformative” will force the Bayes factor to favor H_0 . This was noted by Bartlett (1957) and is sometimes called “Bartlett’s paradox.” As Jeffreys recognized, to avoid this difficulty, priors on parameters being tested (ψ in the discussion here) generally must be proper and not have too big a spread; thus the standard improper priors that he used for estimation are not applicable to testing. We return to the use of improper priors on nuisance parameters in Section 5.3.

5.2 Sensitivity Analysis

Because Bayes factors can be sensitive to the prior, it is important to evaluate the Bayes factor over a range of possibilities. This involves specifying classes of priors to use under H_1 and H_2 , and it also makes the issue of computation more urgent, because many multidimensional integrals [as in Eq. (2)] must be calculated. We mention here several classes of priors that may be used. When there is enough information to yield initial priors with given hyperparameters (such as a $N(\nu, \phi^2)$ prior in which the hyperparameters are ν and ϕ), the hyperparameters may be perturbed (e.g., by halving and doubling ϕ or changing ν to $\nu \pm \phi$) and the Bayes factor recomputed, as in Applications 1, 4, and 5 (see also McCulloch and Rossi 1991).

An important computational device is to use Equation (5) as an approximation to (2) substituted in (1). From this, if we change from an original pair of priors under H_1 and H_2 to a new pair and thus compute a new Bayes factor $B_{12}^{(NEW)}$, we obtain

$$B_{12}^{(NEW)} \doteq B_{12} \cdot \left(\frac{\pi^{(NEW)}(\hat{\theta}_1 | H_1)}{\pi^{(NEW)}(\hat{\theta}_2 | H_2)} \right) \cdot \left(\frac{\pi(\hat{\theta}_2 | H_2)}{\pi(\hat{\theta}_1 | H_1)} \right), \quad (14)$$

with an error of order $O(n^{-1})$, which also holds if (5) is used to evaluate the original Bayes factor B_{12} . Because the ratio of prior ordinates at the MLE is easy to evaluate, (14) makes results for large numbers of priors easy to obtain; see Application 5. The accuracy of the method was explored in detail for the case of testing equality of two binomial proportions by Kass and Vaidyanathan (1992).

The maximum of B_{10} (and thus the maximal evidence against H_0) over classes of priors was discussed by Edwards, Lindman, and Savage (1963) in contrasting Bayes factors with P values and was further developed by Berger and Delampady (1987), Raftery (1988b), and others (see Sec. 8.2). Suppose, first, that ψ is one-dimensional and that there is no additional parameter β . In this case, under H_0 we have $\psi = \psi_0$, and there is no prior. Under H_1 , taking $\psi \sim N(\psi_0, \phi^2)$ leaves only the parameter ϕ^2 to be determined; sensitivity in the Bayes factor is then reflected in its sensitivity to the choice of ϕ . Edwards et al. (1963) computed the maximum of B_{10} over all choices of ϕ , which we will denote by $B_{10}^{(NORMMAX)}$. There is no corresponding minimum; it is generally the case that $B_{10} \rightarrow 0$ as $\phi \rightarrow \infty$.

Edwards et al. (1963) also considered the maximum of B_{10} over all possible priors on ψ , which in simple cases occurs for the prior with all its mass at the MLE $\hat{\psi}$. This class is wider, but it can be too large and can yield bounds on B_{10} that are too big. As a result, Berger and Delampady (1987) examined the class of “symmetric unimodal” priors, which are symmetric about ψ_0 and nondecreasing in $|\psi - \psi_0|$, and obtained an expression for the maximum of B_{10} . They also considered other classes, which opens up the question of what class should be used. Note that $B_{10}^{(NORMMAX)}$ is also the maximum of B_{10} over all priors that are scale mixtures of the $N(\psi_0, \phi^2)$ distribution. This seems a reasonably large and interesting class. Berger and Delampady (1987) showed that $B_{10}^{(NORMMAX)}$ is not very different from the bound obtained using the symmetric unimodal class. A computationally simple approximation to $B_{10}^{(NORMMAX)}$ may be obtained using Laplace’s method [Eq. (5)] together with the argument leading to equation (15); the formula may be found in Kass and Raftery (1993).

When ψ may be a vector and there is also, under both H_0 and H_1 , a nuisance parameter β (again possibly a vector), the problem is more complicated. Jeffreys (1961, pp. 249–250) showed that under certain conditions, the prior on the nuisance parameter β is much less relevant than that on the parameter being tested ψ if the simplification (13) is used together with a priori independence under H_1 . Kass and Vaidyanathan (1992) extended and sharpened the result, which we now describe, by considering local alternatives (ψ near ψ_0).

In addition to the simplification of (13) together with a priori independence of β and ψ under H_1 , assume that β and ψ are *null-orthogonal* in the sense that the Fisher information matrix $J(\beta, \psi)$ is block-diagonal when $\psi = \psi_0$. Kass and Vaidyanathan (1992) noted that this condition often holds, at least approximately, and by a transformation of parameters it can always be made valid. They assumed that $\hat{\psi} - \psi_0 = O(n^{-1/2})$, as would be the case if the “true” value of ψ were either ψ_0 or a neighboring alternative ψ_n such that $\psi_n - \psi_0 = O(n^{-1/2})$: when this situation does not hold, the Bayes factor will quickly become decisive, and issues of asymptotic approximation will no longer be of concern. Now let B_{10} be the Bayes factor for one prior on β and let B_{10}^* be the Bayes factor for a different prior. Then (under suitable regularity conditions),

$$B_{10} = B_{10}^* \cdot (1 + O(n^{-1})). \quad (15)$$

That is, up to order $O(n^{-1})$, the Bayes factor no longer depends on the choice of prior on β . Thus when the parameters are null-orthogonal, sensitivity analysis may be confined to examination of priors on ψ . In practice if, under H_1 , β and ψ are approximately uncorrelated (they are “observed-orthogonal,” which may be easier to check), then the relative insensitivity to choice of prior on β may be expected to hold (this was done in Application 5). Furthermore, the same argument (under the null, or for local alternatives) shows that (15) holds to order $O(n^{-1/2})$ even if the parameters are not null-orthogonal.

When there is little prior information, Raftery (1993c) argued that subjective priors will often be relatively flat in the region where the likelihood is large and that their impact on Bayes factors for the comparison of both nested and non-nested models should be small. Thus an alternative to careful and time-consuming prior elicitation in this case may be to use a set of priors constructed to have this property. Raftery (1993c) showed how this can be done for generalized linear models. These models relate a dependent variable y_i to independent variables $x_i = (x_{i1}, \dots, x_{iJ})$, where $x_{i1} = 1$ in such a way that $E[y_i | x_i] = \mu_i$, $\text{var}[y_i | x_i] = \sigma^2 v(\mu_i)$, and $g(\mu_i) = x_i \beta$, where $\beta = (\beta_1, \dots, \beta_J)^T$. The idea is to narrow down an initial class of baseline priors for β , here taken to be such that when $g(\mu) = \mu$, $v(\mu) = 1$, and the variables have been standardized to have mean zero and variance 1, the β_j are independent normal a priori with $\beta_1 \sim N(\nu, \eta^2)$ and $\beta_j \sim N(0, \phi^2)$ ($j = 2, \dots, J$). Bayes factors tend to be insensitive to ν and η (as indicated by the remarks following (15)), but they can be quite sensitive to the choice of ϕ .

One way of defining a reasonable range of priors is to require that the ratio of prior ordinates at the MLE, given in Equation (14), not be too far from 1 for any of the possible values of $\hat{\beta}$, namely $|\hat{\beta}_2| \leq 1$ when $J = 2$. We wish the same to be true when H_0 and H_1 are not nested. This corresponds to requiring that the prior not contribute much evidence in favor of either model, whether the models being compared are nested or not. This requirement involves a tradeoff: for nonnested models, it implies that ϕ be large, whereas for nested models it implies that ϕ not be too large. Balancing these two desiderata in a certain sense gives $\phi = e^{1/2} = 1.65$, and requiring that the priors not contribute evidence “worth more than a bare mention” beyond what is unavoidable leads to the range $1 \leq \phi \leq 5$. The resulting priors are then transformed back to the original scale for the variables; results for other choices of $g(\mu)$ and $v(\mu)$ are obtained by weighting the cases appropriately.

The result is what Raftery (1993c) calls a *reference set of proper priors* for generalized linear models. These are used in the GLIB software; see Section 9. Although they are mildly data-dependent, they do have properties that one would associate with genuine subjective data-independent priors that represent a small amount of prior information. Similar reasoning can be applied to other classes of baseline priors and to other models.

5.3 Bayes Factors with Improper Priors

In Section 5.1 we indicated that improper priors on parameters of interest (ψ when we have $H_0: \psi = \psi_0$) are prob-

lematic because, when used under H_1 and not under H_0 , they force B_{10} to become zero. But Jeffreys (1961) used improper priors on nuisance parameters appearing in both null and alternative models (e.g., in testing the value of a Normal mean $H_0: \mu = \mu_0$, he took the prior on μ under H_1 to be proper but set $\pi(\sigma) = 1/\sigma$; p. 268). This leads to an improper predictive distribution specified by (2), but the value of (2) for the given data remains well defined, so this impropriety did not seem to bother Jeffreys; others also are untroubled by the procedure (Moulton 1991; Robert 1992; Robert and Caron 1992). Equation (15) shows that in many cases the choice of prior on the nuisance parameter does not greatly affect the results.

Some authors have used improper priors for all parameters appearing in the models. This has the problem that flat priors are specified only up to an undefined multiplicative constant. Thus the Bayes factor in this case also contains undefined constants. One effort to resolve this difficulty is the “imaginary training sample device” of Smith and Spiegelhalter (1981) and Spiegelhalter and Smith (1982). This consists of imagining that a data set is available that involves the smallest possible sample size permitting a comparison of H_0 and H_1 and provides maximum possible support for H_0 , and then arguing that $B_{10} = (1 + \varepsilon)^{-1}$, where $\varepsilon \geq 0$ is small; they took $\varepsilon = 0$. This yields a value for the ratio of constants.

The authors of several published applications of the method found it useful (Akman and Raftery 1986a; Racine, Grieve, Fluhler, and Smith 1986; Raftery 1987, 1988; Raftery and Akman 1986b). The dimensionalities of the alternative models in these examples were not very different, however, and it is not clear how the method will perform in the more difficult case where this is not so.

Another solution is to set aside part of the data to use as a training sample which is combined with the improper prior distribution to produce a prior prior distribution. The Bayes factor is then computed from the remainder of the data. This idea was introduced by Lempers (1971), and other implementations have been suggested more recently under the names partial Bayes factors (O’Hagan 1991), intrinsic Bayes factors (Berger and Perrichi 1993), and fractional Bayes factors (O’Hagan 1995).

6. ACCOUNTING FOR MODEL UNCERTAINTY

Practical model-building often involves far more than the comparison of two models; there are usually many other choices to be made. For example, in regression the analyst must choose the independent variables, decide which if any of the observations are outliers, and determine how if at all to transform the variables. Each possible combination of choices defines a different model, so that the model-building process consists of comparing many competing models. Strategies for doing this are commonly guided by a series of significance tests, often based on the approximate asymptotic distribution of a test statistic.

There are several problems with this. The sampling properties of the overall strategy, as distinct from those of the individual tests, are not well understood (Freedman 1983; Miller 1984, 1990). The models being compared are often not nested. Power considerations are usually not taken into

account when setting significance levels; indeed, the power characteristics of the tests are often unknown. Any approach that selects a single model and then makes inference conditionally on that model ignores the uncertainty involved in model selection, which can be a big part of overall uncertainty. This leads to underestimation of the uncertainty about quantities of interest, sometimes to a dramatic extent. See Application 4 for an example of this.

All of these difficulties can be avoided, at least in principle, if one adopts a Bayesian approach and calculates the posterior probabilities of all the competing models, which follow directly from the Bayes factors (e.g., Leamer 1978; Stewart 1987). A composite inference can then be made that takes account of model uncertainty in a simple and formally justifiable way. In Section 6.1 we review this approach, and in Sections 6.2–6.5 we discuss methods for handling the difficult situation where the number of models is very large.

6.1 Basic Ideas

When several models are being considered, Bayes factors yield their posterior probabilities as follows. Suppose that $(K + 1)$ models, H_0, H_1, \dots, H_K , are being considered. Each of H_1, \dots, H_K is compared in turn with H_0 , yielding Bayes factors B_{10}, \dots, B_{K0} . Then the posterior probability of H_k is

$$\text{pr}(H_k|\mathbf{D}) = \alpha_k B_{k0} / \sum_{r=0}^K \alpha_r B_{r0}, \quad (16)$$

where $\alpha_k = \text{pr}(H_k)/\text{pr}(H_0)$ is the prior odds for H_k against H_0 ($k = 0, \dots, K$); here $B_{00} = \alpha_0 = 1$. Taking all the prior odds α_k equal to 1 is a natural choice, but other values of α_k may be used to reflect prior information about the relative plausibility of competing models.

The posterior model probabilities given by Equation (16) lead directly to solutions of the prediction, decision-making and inference problems that take account of model uncertainty. For a quantity of interest Δ that is well-defined for every model, the posterior density given model H_k is $\text{pr}(\Delta|\mathbf{D}, H_k) = \int \text{pr}(\Delta|\mathbf{D}, \theta_k, H_k) \text{pr}(\theta_k|\mathbf{D}, H_k) d\theta_k$. This can be used to make inferences about Δ conditionally on model H_k , but instead we may use the posterior density of Δ without conditioning, namely

$$\text{pr}(\Delta|\mathbf{D}) = \sum_{k=0}^K \text{pr}(\Delta|\mathbf{D}, H_k) \text{pr}(H_k|\mathbf{D}) \quad (17)$$

(Leamer, 1978, p. 117). This accounts for the uncertainty about model form by weighting the conditional posterior densities according to the posterior probabilities of each model. The posterior mean and standard deviation of Δ are as follows (Leamer, 1978, p. 118):

$$E[\Delta|\mathbf{D}] = \sum_{k=0}^K E[\Delta|\mathbf{D}, H_k] \cdot \text{pr}(H_k|\mathbf{D}) \quad (18)$$

and

$$\begin{aligned} \text{var}[\Delta|\mathbf{D}] &= \sum_{k=0}^K (\text{var}[\Delta|\mathbf{D}, H_k] \\ &+ (E[\Delta|\mathbf{D}, H_k])^2) \cdot \text{pr}(H_k|\mathbf{D}) - E[\Delta|\mathbf{D}]^2. \end{aligned} \quad (19)$$

Racine et al. (1986) showed how this method can be used to make inference about a treatment effect in the presence of uncertainty about the existence of a carryover effect.

The decision-making problem is solved by maximizing the posterior expected utility of each course of action considered. The latter is equal to a weighted average of the posterior expected utilities conditional on each of the models, with the weights equal to the posterior model probabilities $\text{pr}(H_k|\mathbf{D})$. Smith (1991) discussed the situation where a model is to be chosen and then a decision made. Our view is that if possible, a single model should not be selected before decision making, and that model uncertainty should be accounted for in the calculation of posterior expected utilities.

Much of the literature on statistical analysis in the presence of a set of rival models has focused on the selection of a single model. Equation (17) shows that selecting a single model and proceeding conditionally on it may be reasonable if one of the $\text{pr}(H_k|\mathbf{D})$ is close to unity or if the sum is dominated by models for which the values of $\text{pr}(\Delta|\mathbf{D}, H_k)$ are similar. If not, then analyses conditional on a single selected model fail to take account fully of uncertainty about structure and so may well underestimate the uncertainty associated with their conclusions. This can lead to, for example, policy choices that are riskier than one thinks (Hodges 1987).

6.2 Occam's Window

Despite the importance of model uncertainty and the existence of a general strategy for dealing with it, at least since the work of Leamer (1978), there have been three major obstacles to the widespread adoption of the method outlined in the previous section. The first is the difficulty of calculating Bayes factors; we have shown in Section 4 that there is now a range of feasible computational strategies for doing this.

The second obstacle is that the number of terms in Equation (17) can be enormous. For example, in regression with n cases and J candidate independent variables, considering all possible subsets, the possibility of outliers and four possible transformations of each variable gives an initial set of around $(2^J \times \binom{n}{O_{\max}} \times 4^{J+1})$ models, where O_{\max} is the maximum number of possible outliers envisaged. Even for a relatively small problem, with $n = 40$, $J = 12$, and $O_{\max} = 5$, this is on the order of 10^{16} models.

The third obstacle is that prior distributions for the parameters must be specified for each model. Various possible ways around this now exist. One approach is to use the Schwarz criterion, relying on the result of Kass and Wasserman (1995) that this gives an accurate approximation for a particular, reasonable prior. Another way is to specify prior distributions for one or several "big" models within which all or most of the models considered are nested, and then obtain the priors for the nested models by conditioning on the constraints that define them (as in Raftery 1993c).

In this section and the next we describe two general algorithmic approaches to solving this problem that lead to feasible methods. The first, known as "Occam's window," was proposed by Madigan and Raftery (1994) and selects a subset of the models initially considered. This involves averaging over a much smaller set of models than in (17),

thereby facilitating effective communication of model uncertainty.

Those authors argue that if a model is far less likely a posteriori than the most likely model, then it has been discredited and should no longer be considered. Thus models not belonging to

$$\mathcal{A}' = \left\{ H_k: \frac{\max_l \{ \text{pr}(H_l | \mathbf{D}) \}}{\text{pr}(H_k | \mathbf{D})} \leq C \right\}$$

should be excluded from Equation (17), where C is a fairly large number ($\gg 1$) chosen by the data analyst, such as $C = 20$. Appealing to Occam's razor, they also exclude from (17) complex models that receive less support from the data than their simpler counterparts, namely those belonging to

$$\mathcal{B} = \left\{ H_k: \exists H_l \in \mathcal{A}, H_l \subset H_k, \frac{\text{pr}(H_l | \mathbf{D})}{\text{pr}(H_k | \mathbf{D})} > 1 \right\}.$$

Then Equation (17) is replaced by

$$\text{pr}(\Delta | \mathbf{D}) = \frac{\sum_{H_k \in \mathcal{A}} \text{pr}(\Delta | H_k, \mathbf{D}) \text{pr}(\mathbf{D} | H_k) \text{pr}(H_k)}{\sum_{H_k \in \mathcal{A}} \text{pr}(\mathbf{D} | H_k) \text{pr}(H_k)},$$

where $\mathcal{A} = \mathcal{A}' \setminus \mathcal{B}$.

One search strategy to identify the models in \mathcal{A} consists of a sequence of pairwise comparisons of nested models. If a model is rejected in favor of a larger one, then all the models nested within it are also rejected. Also, if there is evidence for H_0 , then H_1 is rejected; but to reject H_0 , we require strong evidence for the larger model, H_1 . If the evidence is inconclusive (falling in "Occam's window"), then neither model is rejected.

Typically the number of terms in (17) is reduced to 25 or less, and often to as few as 1 or 2. This procedure mimics the evolutionary process of model selection that is typical of science. The final solution is fairly independent of the initial class considered, in the sense that most initial classes that contain \mathcal{A} give the same result.

6.3 Markov Chain Monte Carlo Model Composition (MC³)

Madigan and York (1992) proposed approximating (17) using MCMC methods, generating a process that moves through model space. They constructed an irreducible Markov chain $\{H(t)\}$, $t = 1, 2, \dots$ with state space \mathcal{H} and equilibrium distribution $\text{pr}(H_i | \mathbf{D})$, where \mathcal{H} is the space of models considered. Then for any function $u(H_i)$, if this Markov chain is simulated for $t = 1, \dots, m$, the average,

$$\hat{U} = \frac{1}{m} \sum_{t=1}^m u(H(t)),$$

converges with probability 1 to $E[u(H)]$ as $m \rightarrow \infty$. To compute (17) in this way, they set $u(H) = \text{pr}(\Delta | H, \mathbf{D})$.

To construct the Markov chain, for each model H they defined a neighborhood $\text{nb}(H)$ consisting of H itself and the models that differ from H by just one parameter. They defined a transition matrix \mathbf{R} by setting $\mathbf{R}(H \rightarrow H') = 0$ for all $H' \notin \text{nb}(H)$ and $\mathbf{R}(H \rightarrow H')$ constant for all $H' \in \text{nb}(H)$. H' is then drawn from $q(H \rightarrow H')$ and accepted with probability

$$\min \left\{ 1, \frac{\text{pr}(H' | \mathbf{D})}{\text{pr}(H | \mathbf{D})} \right\}.$$

Otherwise, the chain stays in state H . Madigan and York (1992) reported that this process is highly mobile, and that runs of 10,000 or less are typically adequate.

George and McCulloch (1993) proposed a similar method in which the chain moves through both model space and parameter space at once. To ensure that the chain is irreducible, they never actually eliminated a parameter from the model, but instead set it close to zero with high probability.

6.4 Model Expansion

Draper (1995) proposed the model expansion method, in which one model is selected initially and then generalized to a set of models that include the initially selected one as a special case but that relax some of the structural assumptions underlying it. Equation (17) is then used for inference about quantities of interest, but restricted to the set of models obtained in the generalization step. In regression, for example, a set of variables and functional forms might be selected initially and a normal distribution assumed for the errors. This might be generalized by embedding the normal error distribution in the symmetric power exponential family (Box and Tiao 1962). Model expansion can be continuous, as in that example, or discrete.

Model expansion is useful for taking account of uncertainty about specific structural assumptions in a model but is not designed to take account of the uncertainty inherent in model building when many models are initially considered, as in variable selection in regression.

6.5 Evaluation of Methods

The efficacy of a modeling strategy can be judged by how well the resulting predictive distributions predict future observations (Dawid 1984). Madigan and Raftery (1994) measured predictive performance by splitting complete data into two subsets, one (typically about 25%–50% of the total) used to calculate model probabilities and the other used as a set of test cases. Predictive performance was then measured using the logarithmic scoring rule of Good (1952); see Section 3.2. This method can be used to assess the performance of any method that generates predictive distributions, Bayesian or not, model based or not, statistical or not. With this scoring rule, Equation (17) is guaranteed in a certain sense to give better predictions on average than those based on any individual selected model (Madigan and Raftery 1994).

Model averaging by the Occam's window and MC³ methods has given consistently and substantially better predictions than those methods based on any one model alone, for several data sets. The differences between individual "good" models were smaller than the gain due to taking model uncertainty into account. The MCMC method had better predictive performance than Occam's window, but at the cost of greater computational expense and less easily interpretable results. This is true for discrete graphical models (Madigan and Raftery 1994), linear regression models (Raftery, Madigan, and Hoeting 1993), and survival analysis (Raftery et al. 1995).

7. APPLICATIONS, REVISITED

7.1 Application 1: *E. coli* Mutagenesis

The raw data for each strain of *E. coli* were two pairs of sample sizes and corresponding proportions (n_{i1}, \hat{p}_{i1}) and (n_{i2}, \hat{p}_{i2}) . Here $i = 1, \dots, 13$, with the 13th strain being the one in question, *uvrE*. The data were assumed to be distributed as binomial proportions and were transformed to the logit scale according to $Y_i = \log[\hat{p}_{i1}(1 - \hat{p}_{i2})/(\hat{p}_{i2}(1 - \hat{p}_{i1}))]$. Then Y_i was assumed to be normally distributed, with σ_i^2 taken to be known and equal to the first-order approximate variance based on binomial sampling, namely $\sigma_i^2 = [n_{i1}\hat{p}_{i1}(1 - \hat{p}_{i1})]^{-1} + [n_{i2}\hat{p}_{i2}(1 - \hat{p}_{i2})]^{-1}$.

We write $\psi_i = \log(p_{i2}/(1 - p_{i2})) - \log(p_{i1}/(1 - p_{i1}))$ so that the null hypothesis is $H_0: \psi_{13} = 0$, and under this hypothesis, $Y_{13} \sim N(0, \sigma_{13}^2)$. For the alternative H_1 , we begin with $Y_{13} \sim N(\psi_{13}, \sigma_{13}^2)$ and then, considering the data from the first 12 strains to be directly relevant, we use them to formulate a prior for ψ_{13} (taking the strains to be exchangeable). This may be done by assuming that $\psi_i \sim N(\mu, \tau^2)$, iid, $i = 1, \dots, 13$, so that $Y_{13} \sim N(\mu, \tau^2 + \sigma_{13}^2)$. The quantities μ and τ may then be estimated from the data on the first 12 strains. This was done by maximum likelihood following Kass and Steffey (1989), who used the transformed data (the *uvrE* strain is number 12 in their table 1). The Bayes factor is then

$$B_{10} = \frac{n(y_{13}; \hat{\mu}, \tau^2 + \sigma_{13}^2)}{n(y_{13}; 0, \sigma_{13}^2)},$$

where $n(x; m, v)$ denotes the normal density with mean m and variance v evaluated at x . The result was $B_{10} = .065$, indicating positive evidence in favor of H_0 .

7.2 Application 2: The Hot Hand

Kass and colleague K. Hsiao—hereafter HK—analyzed the data for Vinnie Johnson. (Details may be found in Hsiao 1994.) Under H_0 , we have $Y_i \stackrel{\text{iid}}{\sim} B(n_i, p)$ ($i = 1, \dots, 380$), there being 380 games in the data set. Under the alternative, $Y_i \stackrel{\text{iid}}{\sim} B(n_i, p_i)$ independently for $i = 1, \dots, 380$, and $p_i \stackrel{\text{iid}}{\sim} \text{beta}(a, b)$. HK began by reparameterizing the beta distribution according to $(a, b) = (\xi/\omega, (1 - \xi)/\omega)$. Then, under the beta-binomial, $E(Y_i/n_i) = \xi$, and the binomial becomes a limiting case of the beta-binomial as $\omega \rightarrow 0$. Thinking of this as a nested model and using the notation of Section 5.1, ψ becomes ω and β becomes ξ , which reduces to p when $\omega = 0$. The additional simplification mentioned there is then to take ξ and ω to be a priori independent, with the distribution on ξ the same as that used on p under the binomial. HK put a uniform(0, 1) prior on ξ and p and then examined the Bayes factor for several values of ω .

To think about the values of ω that should be of interest, it is helpful to convert ω to the standard deviation $\text{SD} = \xi(1 - \xi)/(1 + \omega^{-1})$, which may be done by taking ξ equal to Johnson's overall shooting rate, $\hat{\xi} = .43$. Under H_1 , Johnson's underlying game shooting ability varies from good games with high values of p_i to bad games with low values. To consider the variation sufficiently large to represent interesting swings in ability, it seems to us that the p_i 's would

have to vary frequently by about $\pm .05$. That is, values of SD of at least .05 (roughly) would be needed to represent a meaningful amount of hot-handedness. With this in mind, here are several values of B_{10} :

ω	SD	B_{10}
.005	.035	.16
.01	.049	.017
.03	.085	3×10^{-7}

From these calculations, the evidence does not rule out small degrees of game-to-game extra-binomial variability, but it is strongly against substantial game-to-game extra-binomial variability (as represented by the beta-binomial). In addition, HK used a normal prior on ω for which "the information in the prior equals the information in one observation" as described in Section 4.1.3 (here the prior was half-normal, because $\omega \geq 0$) and obtained $B_{10} = \frac{1}{19}$ using Laplace's method; the Schwarz criterion gave $B_{10} \approx \frac{1}{21}$. (A more accurate Laplace approximation accounting for $\hat{\omega}$ being at the boundary of the parameter space is described in Hsiao 1994; it yielded $B_{10} = \frac{1}{62}$.)

The uniform prior on ξ and p is a natural choice, but B_{10} may be sensitive to it. HK showed that this sensitivity is not enough to call the conclusions into question; indeed, the parameters are null-orthogonal, as discussed in Section 5.2.

Thus these data show decisively that Vinnie Johnson did not have a "hot hand" in the sense that there was no substantial game-to-game extra-binomial variability in his shooting percentage. (The data may be obtained by sending the e-mail message "send Vinnie.Johnson from data" to statlib@stat.cmu.edu.)

7.3 Application 3: Ozone Exceedances

Success of the regulatory efforts of the Texas Air Control Board would be indicated by a decrease in the rate of occurrence of exceedances. If there were no trend, then the data would be close to being from a homogeneous Poisson process; we denote this model by H_0 . An alternative hypothesis is that the exceedance rate has been decreasing smoothly and gradually. This may be represented by the log-linear Poisson process, $H_1: \lambda(s) = \rho e^{-\gamma s}$, where $\lambda(s)$ is the rate of occurrence at time s and $\gamma > 0$. Another possibility is that the exceedance rate decreased fairly abruptly within a short time period. This may be represented by the change point Poisson process, $H_2: \lambda(s) = \lambda_1$ if $0 \leq s \leq \tau$ and $\lambda(s) = \lambda_2$ if $\tau < s \leq T$.

Raftery (1989) calculated Bayes factors for these hypotheses using improper reference priors and the imaginary training sample device of Spiegelhalter and Smith (1982). This is probably a reasonable approximation to vague prior information here, because degrees of freedom involved in the comparisons are small. The Bayes factor for a gradual change against no trend, B_{10} , is .02, indicating strong evidence (odds of 50) *against* a smooth decrease; the Bayes factor corresponds to strong evidence *for* the null hypothesis. In contrast, the Bayes factor for an abrupt change against no trend, B_{20} , is 2.75, indicating evidence for an abrupt change that is "worth no more than a bare mention." Note

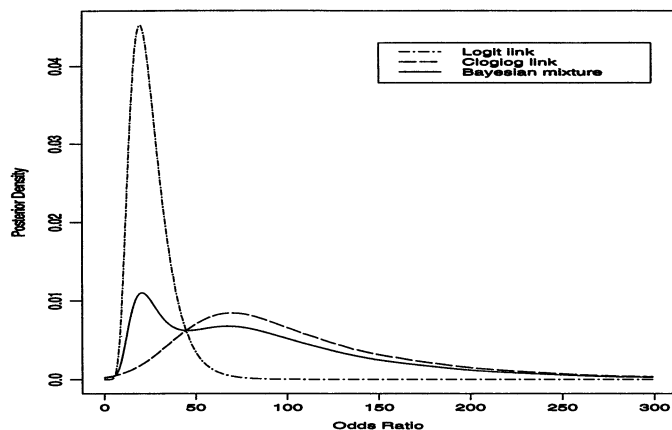


Figure 1. Posterior Distributions of the Odds Ratio in the Educational Transitions Application: The results are indistinguishable for all the prior distributions (i.e., values of ϕ) considered.

that if there was a change, it is much more likely to be abrupt than smooth by odds of $B_{21} = B_{20}/B_{10} = 135$. This is a nonnested model comparison, and one can see how easy it is with the Bayesian method.

The finding of some evidence for an abrupt decrease in the exceedance rate gave rise to the suggestion that this was due not to a change in the underlying ozone levels, but rather to a change in measurement technology that led to lower variances and hence to fewer extreme values. Subsequent exploratory analyses of Smith (1989, rejoinder) suggested that this was plausible, especially because there were indeed changes in measurement instruments that could have led to such a change (Fairley 1989). It would be possible to expand the model to take this into account, and then to test for it explicitly by calculating the Bayes factors for the expanded model against the previously considered models. This example shows that Bayes factors are not restricted to tests of previously formulated hypotheses, but can also be used to guide an empirical model-building process.

In this example Bayes factors have the advantage over frequentist approaches of being *technically simpler*. The comparisons between H_0 and H_1 , and between H_0 and H_2 both lead to nonstandard frequentist testing problems, even though their Bayesian solution is straightforward. Under H_1 , the expected total number of events over all time is finite, so that the usual asymptotic arguments do not apply. Under H_2 , the likelihood is highly discontinuous as a function of the change point τ , so that the usual asymptotics based on smoothness of the likelihood do not apply at all. To see how much more cumbersome the frequentist approach to testing H_0 against H_2 is, compare the frequentist analyses of Worsley (1986) and Akman and Raftery (1986b) with the Bayesian analysis of Raftery and Akman (1986). And we know of no non-Bayesian testing procedures for comparing the non-nested and “irregular” models H_1 and H_2 . In particular, the methods of Cox (1961, 1962) do not apply, because the regularity conditions that they require do not hold.

7.4 Application 4: Educational Transitions

We address the question of interest in this application using part of the longitudinal data collected by Greaney and

Kelleghan (1984). In 1967, they selected a random sample of 11-year-old children in Irish elementary schools and followed them through the rest of their educational careers. Of the 441 students who entered second-level education, 230 completed it by taking the Leaving Certificate Examination and 211 did not. Social class background and ability were both measured prospectively at age 11 by continuous-valued variables, whereas type of school was a dichotomous variable, indicating whether a secondary or vocational school was attended.

We use logistic regression with the reference set of proper priors described in Section 5.2. It is clear that social class background does have an effect, with a Bayes factor exceeding 2,000 for the model that includes it and ability against the model that includes only ability. Thus Greaney and Kelleghan’s conclusion about meritocracy was unwarranted. Also, gender has no effect. A frequentist analysis would say merely that the gender effect is “not significant,” but the Bayesian analysis provides strong evidence (odds of more than 20) for the simpler model from which the gender effect is absent; as in Applications 1 and 2, the Bayes factor provides evidence for a null hypothesis.

The main uncertainty is about the link function. The data provide some evidence (a Bayes factor of 3.7) for the complementary log-log (hereafter termed cloglog) link against the logit link, but this is not strong. This Bayes factor varied by less than .02 over the entire range of priors considered. The models corresponding to different links are not nested, so that a formal non-Bayesian comparison between these models would be cumbersome while the Bayesian comparison is straightforward.

A question of policy interest is the difference between the attainments of secondary and vocational school students, after controlling for the other variables. Here we consider the ratio of the odds on completion of second-level education for secondary school students to that for vocational school students when both social class and ability are high (71 and 127, corresponding to the 95th percentile in each case). We will refer to this quantity simply as the odds ratio. Even though the data do not distinguish clearly between the two link functions, the latter imply very different things about the odds ratio. With the logit link function, the MLE of the odds ratio is 23; with the cloglog link, the MLE is 91—about four times larger.

Approximate posterior distributions of the odds ratio under the two models are shown in Figure 1, as is the combined inference using the Bayesian mixture of Equation (17); the corresponding quantiles are shown in Table 1. The posterior distributions under the two models are very different; the

Table 1. Posterior Quantiles of the Odds Ratio in the Educational Transitions Application

Model	Quantile		
	.025	.5	.975
logit link (H_3)	10	23	52
cloglog link (H_4)	24	94	338
Bayesian mixture	13	77	314

distribution for the logit link is more concentrated and centered about a lower value. The combined posterior distribution has both the peak around 20 from the logit link and the long tail from the cloglog link. The combined 95% interval in Table 1 is close to the union of the intervals corresponding to the individual models but is somewhat shorter; this is intuitively reasonable.

For other analyses of these and related data, see the work of Raftery and Hout (1985, 1993). The data may be obtained by sending the e-mail message "send irish.ed from data" to statlib@stat.cmu.edu.

7.5 Application 5: Predicting Working Memory Failure

The analysis, including prior elicitation, was reported in detail by Carlin et al. (1992). Here we discuss several important issues: the computational burden and the effectiveness of Laplace's method, the difficulty and effectiveness of eliciting the prior and assessing how sensitive the results are to it, the use of influence diagnostics, and the role of the Schwarz criterion.

Computing (2) involved evaluating a one-dimensional integral nested within a five-dimensional integral for the number-of-conditions model and a six-dimensional integral for the query-complexity model. Laplace's method was programmed in the *S* language; the result was $\hat{B}_{21} = 8.0$ in favor of the query-complexity model. This was checked with subregion-adaptive integration, which yielded $B_{21} = 10.8$. Thus Laplace's method was somewhat inaccurate but clearly accurate enough for the inferential purpose at hand. (Subregion-adaptive integration was more difficult and was done by Alan Genz.)

Simple forms were used for the priors, but there were still 15 hyperparameters to be determined. Some of the information for this came from another related experiment, but some was based on the experience of the investigators and was really rough guesswork. Sensitivity analysis was carried out by shifting all prior means by one prior standard deviation unit in each direction and then doubling and halving each prior standard deviation, which resulted in approximately 10^7 alternative prior distributions. Even for this large number of priors, computations were performed quite easily using the Laplace approximation method of Equation (5). After eliminating certain of the alternative priors as unreasonable, the minimal value obtained over all remaining priors was $\hat{B}_{21} = 3.0$.

A concern was the influence that individual subjects might have on the results (via the normal distributional assumption of subject ability effects). Leave-one-out influence diagnostics were computed using the Laplace approximation methods, which are again quite easy to obtain in *S*. A small number of subjects in the experiment did carry much of the comparative information. Taking these sensitivity analyses into account, the overall conclusion was that there was "some evidence, though not strong evidence" in favor of the query complexity model. The Schwarz criterion was $S = 3.6$, giving $B_{21} \approx 6$, which was consistent with the overall conclusion and thus could be used as a summary.

8. ISSUES AND CONTROVERSIES

8.1 Why Test Sharp Hypotheses?

We introduced Bayes factors as a way of assessing the evidence in favor of a scientific theory. Our statement that Jeffreys's approach computes "the posterior probability that one of the theories is correct" invites argument, however. Some would say that theories are never correct, and thus any approach that assumes they are must be flawed. Those who make this argument generally prefer the use of interval estimates.

Certainly we would agree with Jeffreys (1961, pp. 389–390) and Edwards et al. (1963, p. 235) when they say that often hypothesis testing is not applicable and estimation is more appropriate. We also recognize a legitimate worry on the part of many statisticians that empirical models are often taken too seriously, and that poor models are sometimes accepted merely because they fit the data better than other models that are even worse.

But some authors see the introduction of sharp hypotheses as "silly" (e.g., Gelfand, Dey, and Chang, 1992). This viewpoint ignores the way scientific investigations usually proceed: Though one rarely believes a scientific law in an absolute sense, it is a great convenience to speak and to act as if laws are valid. When one says that a certain theory is correct, one means that deviations from it are sufficiently minor to be irrelevant for all practical purposes at hand. Thus theories do become "accepted" for a period, during which they are used to make predictions about new phenomena. Jeffreys (1961, p. 391) pointed out that the best available law gets used for future calculations even when discrepancies are found, noting that "there has not been a single date in the history of the law of gravitation when a modern significance test would not have rejected all laws and left us with no law."

The admission that minor discrepancies may exist between theory and data does not imply that estimation is more natural than testing. For there are many ways in which a theory might err, and in representing the theory by a statistical model one should not presume that possible errors are eliminated simply by letting some parameter ψ take values other than a particular ψ_0 . Simple laws are preferred partly because one cannot be sure that by including more parameters to model one kind of error, the law will succeed better in some unanticipated new situation. In some cases it may be worthwhile to include additional parameters, but this is an empirical question. In fact, it is exactly this empirical question that Bayes factors are supposed to answer. They do so by comparing predictions made by the simpler and the more complicated theories.

For example, in Application 1 an alternative to testing equality of the binomial proportions would be to assume that they are different and to ask how different they are. But how would this correspond to the scientific theory? One would have to say something like "we estimate the difference between the log odds to be less than .38 (with probability .95), which indicates that deviations from the error-prone DNA repair mechanism are quite small." But this is merely an elaboration of the more succinct "evidence in favor of the error-prone DNA repair mechanism." Careful exami-

Table 2. Approximate Minimum t Values for Different Grades of Evidence and Sample Sizes

Evidence for H_1	$2 \log B_{10}$ (Approximation (21))	$n =$			
		30	100	1,000	10,000
Positive	0–6	1.84	2.15	2.63	3.03
Strong	6–10	3.07	3.20	3.59	3.90
Decisive	>10	3.66	3.82	4.11	4.38

nation of the data using either approach ought to lead an analyst to appropriate conclusions. We find that expressing results in terms of Bayes factors is simpler. It is also reassuring that, as Berger and Delampady (1987) noted, in testing a normal mean the point null is a good approximation to an interval null as long as the width of the interval is less than about one-half of the standard error of the sample mean (see also Dickey 1976 and Verdinelli and Wasserman 1993a). (For related discussion along these lines see Jeffreys 1961, p. 367, Raftery 1992, and Zellner 1987.)

It has been argued that a comprehensive account of Bayesian model selection requires decision theory (Kadane and Dickey 1980; Smith 1991). The approach discussed here avoids the introduction of utilities, which would bring with it another layer of sensitivity concerns.

8.2 Bayes Factors Versus Non-Bayesian Significance Testing

There is a substantial literature on the controversy between Bayesian and non-Bayesian testing procedures. This is not a central theme of our article, but we do wish to briefly mention several points that have been made in the literature.

1. There is no reason to expect a P value to be similar to the posterior probability that the null hypothesis is correct. But partly because this misinterpretation of P values is common among nonstatisticians, it is of some interest to compare results. This was done by Jeffreys (1961, pp. 434–435) and many subsequent workers. There is a general feeling that Bayes factors are more conservative than P values, mainly because when comparisons are made, it becomes clear that a P value of .05 cannot represent much evidence against the null (Berger and Mortera 1991 and references therein; Edwards et al. 1963).

2. Frequentist tests tend to reject null hypotheses almost systematically in very large samples, whereas Bayes factors do not. This has been a real problem in sociology, where data sets frequently have thousands of cases. A dramatic example with $n = 113,556$ was discussed by Raftery (1986b, 1995b). There a substantively meaningful model that explained 99.7% of the deviance was rejected by a standard chi-squared test with a P value of about 10^{-120} but was nevertheless favored by the Bayes factor. Faced with this problem, sociologists had taken to ignoring significance tests and using other criteria of reasonableness and common sense in comparing models (e.g., Fienberg and Mason 1979; Grusky and Hauser 1984). Bayes factors are now widely used in sociology, usually with BIC as an approximation.

3. Bayes factors, like Bayesian procedures generally, follow the likelihood principle (e.g., Berger and Wolpert 1984). As

a result, in settings such as clinical trials where cases may accrue sequentially, Bayes factors may be applied without concerns about unscheduled analysis of the data (see, for example, Berger and Berry 1988 and Cornfield 1966a, b).

4. Bayes factors can be applied as easily to nonnested models as to nested ones. In contrast, the application of non-Bayesian significance tests to nonnested models is difficult. The approach of Cox (1961, 1962), which has spawned a large literature, tends to be cumbersome to implement and requires the often arbitrary designation of one of the two nonnested models as the null hypothesis. One way around this arbitrariness is to carry out two tests rather than one test, with each model in turn as the null hypothesis. But there is no guarantee of getting the standard kind of result of a test, namely rejection of one model and nonrejection of the other. Both models may fail to be rejected, in which case it is not clear how to make inferences about quantities of interest, especially if the two models lead to different conclusions. Both models may be rejected (as often happens with large samples), in which case the tests do not provide a comparison between the two models.

5. Non-Bayesian significance tests were developed for the comparison of two models, but practical data analysis often involves far more than two models, at least implicitly. In this case, carrying out multiple frequentist tests to guide a search for the best model can give very misleading results (e.g., Freedman 1983). By allowing us to take into account model uncertainty, Bayes factors can avoid this problem (e.g., Raftery et al. 1993).

With regard to Bayesian calibration of frequentist methods, for large samples the Schwarz criterion may be used to obtain the required value of an approximate t statistic for it to represent strong or decisive evidence. Equation (9) implies that

$$2 \log B_{10} \approx \Lambda - (d_1 - d_0) \log(n). \quad (20)$$

Now if H_0 and H_1 are nested and differ by just 1 degree of freedom, so that $(d_1 - d_0) = 1$, then $\Lambda \approx t^2$ approximately in large samples, where t is the t statistic for (frequentist) testing of H_0 against H_1 , obtained, for example, by dividing $(\hat{\psi} - \psi_0)$ by its large-sample approximate standard error. Thus in this case Equation (20) implies that

$$2 \log B_{10} \approx t^2 - \log(n). \quad (21)$$

The approximate t values corresponding to different grades of evidence (on the cruder scale in Sec. 3.2), and different sample sizes are shown in Table 2 (as in Raftery 1993b; see also Berger 1985). For “positive” evidence, this is $t = \sqrt{\log n}$; for “strong” evidence, it is $t = \sqrt{\log n + 6}$; and for “decisive” evidence, it is $t = \sqrt{\log n + 10}$. Note that the

critical values quickly become larger than the usual cutoffs based on P values as n becomes large.

Atkinson (1978) has noted some instances of Bayes factors favoring the simpler model H_0 even when a more complex model H_1 is correct. But Smith and Spiegelhalter (1980, p. 216) showed that this occurs only when the two models are so close that there is nothing to be lost for predictive purposes by cutting back to the simpler model, so that the Bayes factor functions as a *fully automatic Occam's razor*.

I. J. Good has for many years advocated a Bayes/non-Bayes compromise, consisting essentially of calculating the Bayes factor and using it as a frequentist test statistic. He reviewed these ideas and listed nearly 50 of his own publications on the topic, spanning 40 years (Good 1992). This is a frequentist testing proposal, and there seems to be no reason why it should escape the difficulties that other frequentist tests have. These arise from the way the test statistic is calibrated and not from the choice of test statistic. But Good has also pointed out that his compromise can be inverted, with P values transformed to yield approximate Bayes factors; this is a sensible and useful point.

Goodness-of-fit tests are different from other frequentist tests, because they do not aim to compare two competing models, but rather to detect departures from a null hypothesis even when no alternative has yet been formulated. One might question the use of "rejecting" a hypothesis if there is nothing to put in its place. Many Bayesians do, however, see a useful role for goodness-of-fit tests (Box 1980; Dempster 1971; Rubin 1984). They can be useful for calibrating the diagnostic checks to which a model is subjected and thus guiding the search for a better model. In Jeffreys's view, a model should not be abandoned until a better one (in the posterior model probability sense) is found.

8.3 Bayes Factors Versus the AIC

Akaike (1973) advocated that, given a class of competing models for a data set, one choose the model that minimizes

$$\text{AIC} = -2(\log \text{maximized likelihood}) + 2(\text{number of parameters}). \quad (22)$$

Two main justifications for the AIC have been advanced. The first, due to Akaike (1973), is based on a predictive argument. Suppose that, given current data and a set of possible models, we want the predictive distribution of a future datum. Then, *if the predictive distribution is conditional on a single model and on its estimated parameters*, the AIC picks the model that gives the best approximation, asymptotically, in the Kullback–Leibler sense. But such a predictive distribution is incorrect, because it does not incorporate the uncertainty about parameter values and model form (Aitchison and Dunsmore 1975). Shibata (1976) and Katz (1981) have shown that the AIC tends to overestimate the number of parameters needed, even asymptotically. Thus if one must ignore both parameter uncertainty and model uncertainty when making predictions, it may be worthwhile to have a model that is too big (Shibata 1976; Stone 1979). Related remarks have been given by Zellner (1978) and Stone (1979).

A related argument is that the AIC picks the correct model asymptotically if the complexity of the true model grows with sample size (Shibata 1980, 1981). Typically this is taken to mean that the model grows in one respect (e.g., the order of an autoregressive model) but remains fixed in all other respects (e.g., normality, linearity). Our experience with large data sets in sociology is that the AIC selects models that are too big even when the sample size is large, including effects that are counterintuitive or not borne out by subsequent research.

The second main justification for the AIC, perhaps best described by Akaike (1983), is Bayesian. He wrote that model comparisons based on the AIC are asymptotically equivalent to those based on Bayes factors. But this is true only if the precision of the prior is comparable to that of the likelihood, but not in the more usual situation where prior information is small relative to the information provided by the data. In the latter more usual situation, the Schwarz criterion indicates that the model with the highest posterior probability is the one that minimizes

$$\text{BIC} = -2(\log \text{maximized likelihood}) + (\log N)(\text{number of parameters}). \quad (23)$$

Comparing Equations (22) and (23) indicates that BIC tends to favor simpler models than those chosen by the AIC criterion. [Akaike (1977) proposed a modification of AIC that was consistent in Normal linear regression models; he called the new criterion BIC and, in fact, in this setting it is asymptotically equivalent to (23).]

Findley (1991) gave some cases in which BIC does not yield consistent model selection but AIC does. But these are cases in which standard asymptotics do not apply and thus the theory in Section 4.1.3 leading to the approximation (9) also does not apply.

Linhart and Zucchini (1986) generalized Akaike's (1973) approach, replacing the Kullback–Leibler distance between true and estimated predictive distributions by any arbitrary distance and replacing the quantity to be predicted (the next data point in Akaike's development) by any quantity of interest. Their work shows, for example, that the "2" by which the number of parameters is multiplied in Equation (22) is arbitrary in that it depends crucially on the choices of distance and quantity of interest, and that other choices can lead to quite different multipliers, such as "4." But their approach is open to the same general criticisms as the AIC. In particular, it provides no way of taking account of model uncertainty and so is somewhat at a loss when several models score almost equally well.

9. BIBLIOGRAPHICAL REMARKS AND ADDITIONAL WORK

The works of Jeffreys (1961), Good (1952, 1983, 1985), Mosteller and Wallace (1964), and Zellner (1971) are important basic sources for applications of Bayes factors. The literature on model selection also contains many papers on Bayes factors, particularly in econometrics (e.g., Leamer 1983; McCulloch and Rossi 1991; Poirier 1985; Rossi 1985, 1988; Schotman and van Dijk 1991).

Several Bayesian methods of selecting models without Bayes factors have been proposed, often using some kind of sample-splitting as in cross-validation (Fornell and Rust 1989; Geisser and Eddy 1979; Gelfand et al. 1992).

Aitkin (1991) described so-called “posterior Bayes factors” in which the marginal likelihood $\text{pr}(\mathbf{D}|H_k)$ defined by Equation (2), which is the prior mean of the likelihood, is replaced in Equation (1) by the posterior mean of the likelihood, namely $\int \text{pr}(\mathbf{D}|\theta_k, H_k)\text{pr}(\theta_k|\mathbf{D}, H_k) d\theta_k$. Despite their name, these are quite different in practice from Bayes factors as defined in this article. Several discussants of Aitkin (1991) pointed out that the procedure has little Bayesian justification, does not have any known frequentist optimality properties, and yields counterintuitive results (e.g., discussions by Dawid, Fearn, Goldstein, Lindley, and Whittaker).

We now give a noncomprehensive list of references in which explicit expressions for Bayes factors in different models appear.

Multinomials: Jeffreys (1961); Good (1967); Dickey and Lientz (1970); Gunel and Dickey (1974); Good and Crook (1987); Albert (1990)

Linear models: Zellner (1971); Dickey (1971); Dickey (1975); Zellner and Siow (1980); Smith and Spiegelhalter (1981); Zellner (1984); Broemeling (1985); Draper and Guttman (1987); Mitchell and Beauchamp (1988); Raftery, Madigan, and Hoeting (1993)

Outliers: Pena and Guttman (1992); Pettit and Young (1990)

Logistic regression and log-linear models: Raftery (1986a, 1988b, 1993c); Stewart (1987); Madigan and Raftery (1994); McCulloch and Rossi (1991); Madigan and York (1992); Kass and Vaidyanathan (1992); also the GLIB software

Survival analysis: Raftery, Madigan, and Volinsky (1995)

Multivariate analysis: Dickey (1971); Dickey (1975); Dayal and Dickey (1976); Smith and Spiegelhalter (1981) [multivariate normal models]; Cooper and Herskovits (1992); Cowell, Dawid, and Spiegelhalter (1993), Madigan and Raftery (1994); Madigan, Raftery, York, Bradshaw, and Almond (1994) [discrete graphical models]; Banfield and Raftery (1993) [cluster analysis]; Raftery (1993b) [structural equation models (LISREL)]

Stochastic processes: Dickey and Lientz (1970), Katz (1981) [Markov chains]; Broemeling (1985) [autoregressive models]; Katz (1981) [Markov chains]; Akman and Raftery (1986a) [point processes]; Raftery and Akman (1986) [change-point Poisson processes]; Raftery (1987, 1988) [software reliability]; Le, Raftery, and Martin (1990) [autoregressive models with outliers]. Polson and Roberts (1994) [diffusion processes]

Deterministic models: Raftery, Givens, and Zeh (1995)

Other models: Berry, Evett, and Pinchin (1992) [assessment of forensic evidence]

Software: There is little general purpose software for Bayes factors, although the BIC approximation can often be easily calculated from the output of standard statistical software. GLIB (generalized linear Bayesian modeling)

is an S-PLUS function that returns accurate Bayes factors and posterior probabilities for generalized linear models, as well as inference about the parameters that takes account of model uncertainty. It can be obtained by sending the e-mail message “send glib from S” to statlib@stat.cmu.edu. BICREG and BIC.LOGIT are S-PLUS functions that do approximate Bayesian model selection and accounting for model uncertainty in linear regression and logistic regression, respectively. They can be obtained by sending the messages “send bicreg from S” and “send bic.logit from S” to statlib@stat.cmu.edu.

10. CONCLUSION

We have reviewed applications and developments of a method introduced by Jeffreys more than 50 years ago. The appeal of the concept was that it provided a simple and satisfying description of the process of accepting new scientific laws as replacements for older ones, which Jeffreys illustrated through the use of many examples in his book *Theory of Probability*.

We have tried to show that Jeffreys’s methodology has worn well with time and that the Bayesian approach to hypothesis testing has evolved to fill a niche in modern computationally intensive statistical practice. It applies to a limited but important class of problems in scientific inference and also to the assessment of uncertainty when many models are considered initially. We wish to emphasize the following points:

- From Jeffreys’ Bayesian viewpoint, the purpose of hypothesis testing is to evaluate the evidence in favor of a scientific theory.
- Bayes factors offer a way of evaluating evidence *in favor of* a null hypothesis.
- Bayes factors provide a way of including other information when assessing the evidence for a hypothesis.
- Bayes factors are very general. In particular, they do not require alternative models to be nested.
- Several techniques are available for computing Bayes factors, including asymptotic approximations that are easy to compute using the output from standard packages that maximize likelihoods.
- In “nonstandard” statistical models that do not satisfy common regularity conditions, it can be technically simpler to calculate Bayes factors than to derive non-Bayesian significance tests.
- The Schwarz criterion (or BIC) gives a rough approximation to the logarithm of the Bayes factor, which is easy to use and does not require evaluation of prior distributions. It is well suited for summarizing results in scientific communication.
- When several models are considered initially, Bayes factors can be used to calculate posterior model probabilities, yielding composite estimates or predictions that take account of model uncertainty.
- Algorithms have been proposed that allow model uncertainty to be taken into account when the class of models initially considered is very large.
- Bayes factors are useful for guiding an evolutionary model-building process.

- It is important, and feasible, to assess the sensitivity of conclusions to the prior distributions used. In our applications we have found our conclusions to be robust to the prior in a qualitative sense, but this is not guaranteed to be the case.

Bayes factors have many of the strengths and limitations of the Bayesian approach generally. The essential strength is their solid logical foundation, which offers great flexibility. Some of the practical advantages of this were noted earlier; another is invariance with respect to stopping rules in clinical trials, mentioned in Section 8.2. Recently, advances in computing and the development of methods that take advantage of additional computational power have greatly extended the usefulness of Bayesian methods. Bayes factors can now be computed for a wide variety of models.

The chief limitations of Bayes factors are their sensitivity to the assumptions in the parametric model and the choice of priors. We have discussed ways of doing sensitivity analysis in Section 5, which are illustrated in Applications 4 and 5. Application 4 illustrates the usefulness of a reference set of proper priors for sensitivity analysis, which has been implemented for generalized linear models in the GLIB software; this idea needs to be extended to other classes of models.

The more situation-specific sensitivity analysis in Application 5 is rather cumbersome. It requires close attention to the details of model and prior. It may be argued that this is appropriate; we should not be cavalier in making inferences that depend on our assumptions. The Schwarz criterion (or BIC) may be used for reporting scientific results with other analysis omitted but serving as background support. The question of how much effort should be made before conclusions are drawn arises in any data analysis problem. This is part of the art of applied statistics, but there is room for research that would help the applied statistician decide whether or not to proceed with the full Bayesian analysis.

[Received March 1993. Revised July 1994.]

REFERENCES

- Aitchison, J., and Dunsmore, I. R. (1975), *Statistical Prediction Analysis*, Cambridge, U.K.: Cambridge University Press.
- Aitkin, M. (1991), "Posterior Bayes Factors" (with discussion), *Journal of the Royal Statistical Society, Ser. B*, 53, 111–142.
- Akaike, H. (1973), "Information Theory and an Extension of the Maximum Likelihood Principle," in *Second International Symposium on Information Theory*, eds. B. N. Petrox and F. Caski, Budapest: Akademiai Kiado, p. 267.
- (1977), "On Entropy Maximization Principle," in *Applications of Statistics*, ed. P. R. Krishnaiah, Amsterdam: North-Holland, pp. 27–41.
- (1983), "Information Measures and Model Selection," *Bulletin of the International Statistical Institute*, 50, 277–290.
- Akman, V. E., and Raftery, A. E. (1986a), "Bayes Factors for Non-Homogeneous Poisson Processes With Vague Prior Information," *Journal of the Royal Statistical Society, Ser. B*, 48, 322–329.
- (1986b), "Asymptotic Analysis of a Change-Point Poisson Process," *The Annals of Statistics*, 14, 1583–1590.
- Albert, J. H. (1990), "A Bayesian Test for a Two-Way Contingency Table Using Independence Priors," *Canadian Journal of Statistics*, 4, 347–363.
- Atkinson, A. C. (1978), "Posterior Probabilities for Choosing a Regression Model," *Biometrika*, 65, 39–48.
- Banfield, J. D., and Raftery, A. E. (1993), "Model-Based Gaussian and Non-Gaussian Clustering," *Biometrics*, 49, 803–821.
- Bartlett, M. S. (1957), "Comment on 'A Statistical Paradox' by D. V. Lindley," *Biometrika*, 44, 533–534.
- Berger, J. O. (1985), *Statistical Decision Theory and Bayesian Analysis*, New York: Springer-Verlag.
- Berger, J. O., and Berry, D. A. (1988), "The Relevance of Stopping Rules in Statistical Inference," in *Statistical Decision Theory and Related Topics IV*, Vol. 1, eds. S. S. Gupta and J. O. Berger, New York: Springer-Verlag, pp. 29–47.
- Berger, J. O., and Delampady, M. (1987), "Testing Precise Hypotheses," *Statistical Science*, 3, 317–352.
- Berger, J. O., and Mortera, J. (1991), "Interpreting the Stars in Precise Hypothesis Testing," *International Statistical Review*, 59, 337–353.
- Berger, J. O., and Perrichi, L. R. (1993), "The Intrinsic Bayes Factor for Model Selection and Prediction," technical report, Purdue University, Dept. of Statistics.
- Berger, J. O., and Sellke, T. (1987), "Testing a Point Null Hypothesis: The Irreconcilability of P Values and Evidence," *Journal of the American Statistical Association*, 82, 112–122.
- Berger, J. O., and Wolpert, R. L. (1984), *The Likelihood Principle*, IMS Monograph Vol. 6, Hayward, CA: Institute of Mathematical Statistics.
- Berry, D. A., Evett, I. W., and Pinchin, R. (1992), "Statistical Inference in Crime Investigations Using Deoxyribonucleic Acid Profiling" (with discussion), *Applied Statistics*, 41, 499–531.
- Box, G. E. P. (1980), "Sampling and Bayes's Inference in Scientific Modeling" (with discussion), *Journal of the Royal Statistical Society, Ser. A*, 143, 383–430.
- Box, G. E. P., and Tiao, G. C. (1962), "A Further Look at Robustness via Bayes's Theorem," *Biometrika*, 49, 419–432.
- Broemeling, L. L. (1985), *Bayesian Inference for Linear Models*, New York: Marcel Dekker.
- Carlin, B. P., and Chib, S. (1993), Bayesian Model Choice via Markov Chain Monte Carlo," Research Report 93-006, University of Minnesota, Division of Biostatistics.
- Carlin, B. P., Kass, R. E., Lerch, J., and Huguenard, B. (1992), "Predicting Working Memory Failure: A Subjective Bayesian Approach to Model Selection," *Journal of the American Statistical Association*, 87, 319–327.
- Carlin, B. P., and Polson, N. G. (1991), "Inference for Nonconjugate Bayesian Models Using the Gibbs Sampler," *Canadian Journal of Statistics*, 19, 399–405.
- Chen, M.-H. (1992), "Importance Weighted Marginal Bayesian Posterior Density Estimation," technical report, Purdue University, Dept. of Statistics.
- Chow, G. C. (1981), "A Comparison of the Information and Posterior Probability Criteria for Model Selection," *Journal of Econometrics*, 16, 21–33.
- Cooper, G., and Herskovits, E. (1992), "A Bayesian Method for the Induction of Probabilistic Networks From Data," *Machine Learning*, 9, 309–347.
- Cornfield, J. (1966a), "A Bayesian Test of Some Classical Hypotheses, With Applications to Sequential Clinical Trials," *Journal of the American Statistical Association*, 61, 577–594.
- (1966b), "Sequential trials, sequential analysis, and the likelihood principle," *The American Statistician*, 20, 18–23.
- Cowell, R. G., Dawid, A. P., and Spiegelhalter, D. J. (1993), "Sequential Model Criticism in Probabilistic Expert Systems," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 15, 209–219.
- Cox, D. R. (1961), "Tests of Separate Families of Hypotheses," in *Proceedings of the Fourth Berkeley Symposium*, 1, 105–123.
- (1962), "Further Results on Tests of Separate Families of Hypotheses," *Journal of the Royal Statistical Society, Ser. B*, 24, 406–424.
- Dawid, A. P. (1984), "Present Position and Potential Developments: Some Personal Views. Statistical Theory. The Prequential Approach" (with discussion), *Journal of the Royal Statistical Society, Ser. A*, 147, 178–292.
- (1992), "Prequential Analysis, Stochastic Complexity and Bayesian Inference," in *Bayesian Statistics 4*, eds. J. M. Bernardo et al., Oxford, U.K.: Oxford Science Publications, pp. 109–125.
- Dayal, H. H., and Dickey, J. M. (1976), "Bayes Factors for Behrens–Fisher Problems," *Sankhya*, 38, 315–328.
- De Bruijn, N. G. (1970), *Asymptotic Methods in Analysis*, Amsterdam: North-Holland.
- DeGroot, M. H. (1970), *Optimal Statistical Decisions*. New York: McGraw-Hill.
- Delampady, M., and Berger, J. O. (1990), "Lower Bounds on Bayes Factors for Multinomial Distribution, With Application to Chi-Squared Tests of Fit," *The Annals of Statistics*, 18, 1295–1316.
- Dempster, A. P. (1971), "Model Searching and Estimation in the Logic of Inference," in *Foundations of Statistical Inference*, eds. V. P. Godambe and D. A. Sprott, Toronto: Holt, Rinehart and Winston.
- Dickey, J. M. (1971), "The Weighted Likelihood Ratio, Linear Hypotheses on Normal Location Parameters," *Annals of Mathematical Statistics*, 42, 204–223.

- (1975), "Bayesian Alternatives to the F-Test and Least Squares Estimate in the Normal Linear Model," in *Studies in Bayesian Econometrics and Statistics*, eds. S. E. Fieberg and A. Zellner, Amsterdam: North-Holland, pp. 515–554.
- (1976), "Approximate Posterior Distributions," *Journal of the American Statistical Association*, 71, 680–689.
- (1980), "Approximate Coherence for Regression Model Inference," in *Bayesian Analysis in Econometrics and Statistics*, ed. A. Zellner, Amsterdam: North-Holland, pp. 334–354.
- Dickey, J. M., and Kadane, J. B. (1980), "Bayesian Decision Theory and the Simplification of Models," in *Evaluation of Econometric Models*, eds. J. Kmenta and J. Ramsey, New York: Academic Press, pp. 245–268.
- Dickey, J. M., and Lientz, B. P. (1970), "The Weighted Likelihood Ratio, Sharp Hypotheses About Chances, the Order of a Markov Chain," *Annals of Mathematical Statistics*, 41, 214–226.
- Draper, D. (1995), "Assessment and Propagation of Model Uncertainty" (with discussion), *Journal of the Royal Statistical Society*, Ser. B, 57, 45–98.
- Draper, N. R., and Guttman, I. (1987), "A Common Model Selection Criterion," *Probability and Bayesian Statistics*, ed. R. Viertl, New York: Plenum, pp. 139–150.
- Edwards, W., Lindman, H., and Savage, L. J. (1963), "Bayesian Statistical Inference for Psychological Research," *Psychological Review*, 70, 193–242.
- Evans, M., and Swartz, T. (1995), "Methods for Approximating Integrals in Statistics With Special Emphasis on Bayesian Integration problems," *Statistical Science*, to appear.
- Evetts, I. W. (1991), "Implementing Bayesian Methods in Forensic Science," paper presented at the Fourth Valencia International Meeting on Bayesian Statistics.
- Fairley, D. (1989), "Comment," *Statistical Science*, 4, 381–383.
- Fienberg, S. E., and Mason, W. M. M. (1979), "Identification and Estimation of Age-Period-Cohort Effects in the Analysis of Discrete Archival Data," *Sociological Methodology 1979* (K. F. Schuessler, ed.), San Francisco: Jossey-Bass, pp. 1–67.
- Findley, D. F. (1991), "Counterexamples to Parsimony and BIC," *Annals of the Institute of Statistical Mathematics*, 43, 505–514.
- Fornell, C., and Rust, R. T. (1989), "Incorporating Prior Theory in Covariance Structure Analysis: A Bayesian Approach," *Biometrika*, 54, 249–259.
- Freedman, D. A. (1983), "A Note on Screening Regression Equations," *The American Statistician*, 37, 152–155.
- Garthwaite, P. H. (1992), "Preposterior Expected Loss as a Scoring Rule for Prior Distributions," technical report, University of Aberdeen, Dept. of Statistics.
- Geisser, S., and Eddy, W. F. (1979), "A Predictive Approach to Model Selection," *Journal of the American Statistical Association*, 74, 153–160.
- Gelfand, A. E., and Dey, D. K. (1994), "Bayesian Model Choice: Asymptotics and Exact Calculations," *Journal of the Royal Statistical Society*, Ser. B, 56, 501–514.
- Gelfand, A. E., Dey, D. K., and Chang, H. (1992), "Model Determination Using Predictive Distributions With Implementations via Sampling-Based Methods," in *Bayesian Statistics 4*, eds. J. M. Bernardo et al., Oxford, U.K.: Oxford University Press, pp. 147–167.
- Genz, A., and Kass, R. E. (1993), "Subregion Adaptive Integration of Functions Having a Dominant Peak," technical report, Carnegie Mellon University, Dept. of Statistics.
- Geweke, J. (1989), "Bayesian Inference in Econometric Models Using Monte Carlo Integration," *Econometrica*, 57, 1317–1340.
- George, E. I., and McCulloch, R. E. (1993), "Variable Selection via Gibbs Sampling," *Journal of the American Statistical Association*, 88, 881–889.
- Gilovich, T., Vallone, R., and Tversky, A. (1985), "The Hot Hand in Basketball: On the Misperception of Random Sequences," *Cognitive Psychology*, 17, 295–314.
- Good, I. J. (1950), *Probability and the Weighing of Evidence*, London: Griffin.
- (1952), "Rational Decisions," *Journal of the Royal Statistical Society*, Ser. B, 14, 107–114.
- (1958), "Significance Tests in Parallel and in Series," *Journal of the American Statistical Association*, 53, 799–813.
- (1967), "A Bayesian Significance Test for Multinomial Distributions," *Journal of the Royal Statistical Society*, Ser. B, 29, 399–431.
- (1983), "Explicativity, Corroboration, and the Relative Odds of Hypotheses (#846)," in *Good Thinking—The Foundations of Probability and Its Applications*, Minneapolis: University of Minnesota Press.
- (1985), "Weight of Evidence: A Brief Survey," in *Bayesian Statistics 2*, eds. J. M. Bernardo, M. H. DeGroot, D. V. Lindley, and A. F. M. Smith, New York: Elsevier, pp. 249–269.
- (1992), "The Bayes/non-Bayes Compromise: A Brief Review," *Journal of the American Statistical Association*, 87, 597–606.
- Good, I. J., and Crook, J. F. (1987), "The Robustness and Sensitivity of the Mixed Dirichlet Bayesian Test for 'Independence' in Contingency Tables," *The Annals of Statistics*, 15, 670–693.
- Greaney, V., and Kelleghan, T. (1984), *Equality of Opportunity in Irish Schools*, Dublin: Educational Company.
- Groenewald, P. C. N., and de Waall, D. J. (1989), "Bayesian Tests for Some Precise Hypotheses on Multi-Normal Means," *Australian Journal of Statistics*, 31, 393–408.
- Grusky, D. B., and Hauser, R. M. (1984), "Comparative Social Mobility Revisited: Models of Convergence and Divergence in 16 Countries," *American Sociological Review*, 49, 19–38.
- Gunel, E., and Dickey, J. (1974), "Bayes Factors for Independence in Contingency Tables," *Biometrika*, 61, 545–557.
- Hammersley, J. M., and Handscomb, D. C. (1964), *Monte Carlo Methods*, London: Chapman and Hall.
- Hartigan, J. A. (1992), "Locally Uniform Prior Distributions," technical report, Yale University, Dept. of Statistics.
- Haughton, D. M. A. (1988), "On the Choice of a Model to Fit Data From an Exponential Family," *The Annals of Statistics*, 16, 342–355.
- Hodges, J. S. (1987), "Uncertainty, Policy Analysis and Statistics," *Statistical Science*, 2, 259–291.
- Hsiao, K. (1994), "Bayesian Tests of Extra-Binomial Variability With Emphasis on the Boundary Case," Ph.D. dissertation, Carnegie Mellon University, Dept. of Statistics.
- Jeffreys, H. (1935), "Some Tests of Significance, Treated by the Theory of Probability," *Proceedings of the Cambridge Philosophy Society*, 31, 203–222.
- (1961), *Theory of Probability* (3rd ed.), Oxford, U.K.: Oxford University Press.
- Kadane, J. B., and Dickey, J. M. (1980), "Bayesian Decision Theory and the Simplification of Models," in *Evaluation of Econometric Models*, eds. J. Kmenta and J. Ramsey, New York: Academic Press, 245–268.
- Kadane, J. B., Dickey, J., Winkler, R., Smith, W., and Peters, S. (1980), "Interactive Elicitation of Opinion for A Normal Linear Model," *Journal of the American Statistical Association*, 75, 845–854.
- Kass, R. E. (1993), "Bayes Factors in Practice," *Statistician*, 42, 551–560.
- Kass, R. E., and Greenhouse, J. B. (1989), Comment on "Investigating Therapies of Potentially Great Benefit: ECMO," by Ware (1989), *Statistical Science*, 4, 310–317.
- Kass, R. E., and Raftery, A. E. (1993), "Bayes Factors and Model Uncertainty," Technical Report 571, Carnegie Mellon University, Dept. of Statistics.
- Kass, R. E., and Steffey, D. L. (1989), "Approximate Bayesian Methods in Conditionally Independent Hierarchical Models (Parametric Empirical Bayes Models)," *Journal of the American Statistical Association*, 84, 717–726.
- Kass, R. E., Tierney, L., and Kadane, J. B. (1990), "The Validity of Posterior Asymptotic Expansions Based on Laplace's Method," in *Bayesian and Likelihood Methods in Statistics and Econometrics*, eds. S. Geisser, J. S. Hodges, S. J. Press, and A. Zellner, New York: North-Holland.
- Kass, R. E., and Vaidyanathan, S. (1992), "Approximate Bayes Factors and Orthogonal Parameters, With Application to Testing Equality of Two Binomial Proportions," *Journal of the Royal Statistical Society*, Ser. B, 54, 129–144.
- Kass, R. E., and Wasserman, L. (1992a), "Improving the Laplace Approximation Using Posterior Simulation," Technical Report 566, Carnegie Mellon University, Dept. of Statistics.
- (1995), "A Reference Bayesian Test for Nested Hypotheses and Its Relationship to the Schwarz Criterion," *Journal of the American Statistical Association*, to appear.
- Katz, R. W. (1981), "On Some Criteria for Estimating the Order of a Markov Chain," *Technometrics*, 23, 243–249.
- Le, N. D., Raftery, A. E., and Martin, R. D. (1990), "Robust Order Selection for Autoregressive Models Using Robust Bayes Factors," Technical Report 194, University of Washington, Dept. of Statistics.
- Leamer, E. E. (1978), *Specification Searches: Ad Hoc Inference With Non-experimental Data*, New York: John Wiley.
- (1983), "Model Choice and Specification Analysis," in *Handbook of Econometrics*, Vol. 1, eds. Z. Griliches and M. D. Intriligator, Amsterdam: North-Holland.
- Leonard, T. (1982), Comment on "A Simple Predictive Density Function," *Journal of the American Statistical Association*, 77, 657–658.
- Lempers, F. B. (1971), *Posterior Probabilities of Alternative Linear Models*, Rotterdam: University Press.

- Lewis, S. M., and Raftery, A. E. (1994), "The Laplace-Metropolis Estimator for Bayes Factors via Posterior Simulation," Technical Report 279, University of Washington, Dept. of Statistics.
- Lindley, D. V. (1957), "A Statistical Paradox," *Biometrika*, 44, 187-192.
- (1961), "The Use of Prior Probability Distributions in Statistical Inference and Decisions," in *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Prob.*, ed. J. Neyman, Berkeley, CA: University of California Press, pp. 453-468.
- (1988), "Statistical Inference Concerning Hardy-Weinberg Equilibrium," in *Bayesian Statistics 3*, eds. J. M. Bernardo, M. H. DeGroot, D. V. Lindley, and A. F. M. Smith, Oxford, U.K.: Oxford University Press, pp. 307-326.
- Linhart, H., and Zucchini, W. (1986), *Model Selection*, New York: John Wiley.
- Madigan, D., and Raftery, A. E. (1994), "Model Selection and Accounting for Model Uncertainty in Graphical Models Using Occam's Window," *Journal of the American Statistical Association*, 89, 1535-1546.
- Madigan, D., Raftery, A. E., York, J. C., Bradshaw, J. M., and Almond, R. G. (1994), "Strategies for Graphical Model Selection," in *Selecting Models From Data: AI and Statistics IV*, eds. P. Cheeseman and R. W. Oldford, New York: Springer-Verlag, pp. 91-100.
- Madigan, D., and York, J. (1992), "Bayesian graphical models for discrete data," Technical Report 259, University of Washington, Dept. of Statistics.
- McCulloch, R. E., and Rossi, P. E. (1991), "A Bayesian Approach to Testing the Arbitrage Pricing Theory," *Journal of Econometrics*, 49, 141-168.
- (1993), "Bayes Factors for Nonlinear Hypotheses and Likelihood Distributions," *Biometrika*, 79, 663-676.
- Meng, X. L., and Wong, W. H. (1993), "Simulating Ratios of Normalizing Constants via a Simple Identity," Technical Report 365, University of Chicago, Dept. of Statistics.
- Miller, A. J. (1984), "Selection of Subsets of Regression Variables" (with discussion), *Journal of the Royal Statistical Society*, Ser. A, 147, 389-425.
- (1990), *Subset Selection in Regression*, London: Chapman and Hall.
- Mitchell, T. J., and Beauchamp, J. J. (1988), "Bayesian Variable Selection in Regression," *Journal of the American Statistical Association*, 83, 1023-1032.
- Mosteller, F., and Wallace, D. L. (1964), *Inference and Disputed Authorship: The Federalist*, Reading, MA: Addison-Wesley. Reprinted in 1985 by Springer-Verlag.
- Moulton, B. R. (1991), "A Bayesian Approach to Regression Selection and Estimation, With Application to a Price Index for Radio Services," *Journal of Econometrics*, 49, 169-193.
- Neal, R. M. (1992), "Probabilistic Inference for Artificial Intelligence Using Monte Carlo Methods Based on Markov Chains," technical report, University of Toronto, Dept. of Computer Science.
- Newton, M. A., and Raftery, A. E. (1994), "Approximate Bayesian Inference by the Weighted Likelihood Bootstrap" (with discussion), *Journal of the Royal Statistical Society*, Ser. B, 56, 3-48.
- O'Hagan, A. (1991), Contribution to the discussion of "Posterior Bayes Factors," *Journal of the Royal Statistical Society*, Ser. B, 53, 137.
- (1995), "Fractional Bayes Factors for Model Comparison" (with discussion), *Journal of the Royal Statistical Society*, Ser. B, 57, .
- Parzen, E. (1974), "Some Recent Advances in Time Series Modeling," *IEEE Transactions on Automatic Control*, 19, 723-730.
- Pena, D., and Guttman, I. (1992), "Comparing Probabilistic Methods for Outlier Detection in Linear Models," technical report, University of Toronto, Dept. of Statistics.
- Pettit, L. I., and Young, K. D. S. (1990), "Measuring the Effect of Observations on Bayes Factors," *Biometrika*, 77, 455-466.
- Poirier, D. J. (1985), "Bayesian Hypothesis Testing in Linear Models With Continuously Induced Conjugate Priors Across Hypotheses," in *Bayesian Statistics 2*, eds. J. M. Bernardo, M. H. DeGroot, D. V. Lindley, and A. F. M. Smith, New York: Elsevier, pp. 711-722.
- Polson, N. G., and Roberts, G. O. (1994), "Bayes Factors for Discrete Observations From Diffusion Processes," *Biometrika*, 81, 11-26.
- Racine, A., Grieve, A. P., Fluhler, H., and Smith, A. F. M. (1986), "Bayesian Methods in Practice: Experiences in the Pharmaceutical Industry" (with discussion), *Applied Statistics*, 35, 93-150.
- Raftery, A. E. (1986a), "A Note on Bayes Factors for Log-Linear Contingency Table Models With Vague Prior Information," *Journal of the Royal Statistical Society*, Ser. B, 48, 249-250.
- (1986b), "Choosing Models for Cross-Classifications," *American Sociological Review*, 51, 145-146.
- (1987), "Inference and Prediction for a General Order Statistic Model With Unknown Population Size," *Journal of the American Statistical Association*, 82, 1163-1168.
- (1988a), "Analysis of a Simple Debugging Model," *Applied Statistics*, 37, 12-22.
- (1988b), "Approximate Bayes Factors for Generalized Linear Models," Technical Report 121, University of Washington, Dept. of Statistics.
- (1989), "Are Ozone Exceedance Rates Decreasing?," *Statistical Science*, 4, 378-381.
- (1992), Discussion of "Model Determination Using Predictive Distributions With Implementations via Sampling-Based Methods," in *Bayesian Statistics 4*, eds. J. M. Bernardo et al., Oxford, U.K.: Oxford University Press, 160-163.
- (1993a), "Discussion of Three Papers on the Gibbs Sampler and Other Markov Chain Monte Carlo Methods," *Journal of the Royal Statistical Society*, Ser. B, 53, 85.
- (1993b), "Bayesian Model Selection in Structural Equation Models," in *Testing Structural Equation Models*, eds. K. A. Bollen and J. S. Long, Beverly Hills, CA: Sage, pp. 163-180.
- (1993c), "Approximate Bayes Factors and Accounting for Model Uncertainty in Generalized Linear Models," Technical Report 255, University of Washington, Dept. of Statistics.
- (1995a), "Hypothesis Testing and Model Selection with Posterior Simulation," in *Practical Markov Chain Monte Carlo*, eds. W. R. Gilks, D. J. Spiegelhalter, and S. Richardson, London: Chapman and Hall.
- (1995b), "Bayesian Model Selection in Social Research (with Discussion)," to appear in *Sociological Methodology 1995* (P. V. Marden, ed.), Cambridge, MA: Blackwells.
- Raftery, A. E., and Akman, V. E. (1986), "Bayesian Analysis of a Poisson Process With a Change-Point," *Biometrika*, 73, 85-89.
- Raftery, A. E., and Banfield, J. D. (1990), "Stopping the Gibbs Sampler, the Use of Morphology and Other Issues in Spatial Statistics," *Annals of the Institute of Statistical Mathematics*, 43, 32-43.
- Raftery, A. E., Givens, G. H., and Zeh, J. E. (1995), "Inference From a Deterministic Population Dynamics Model About Bowhead Whale, *Balaena mysticetus*, Replacement Yield" (with discussion), *Journal of the American Statistical Association*, 90, 402-416.
- Raftery, A. E., and Hout, M. (1985), "Does Irish Education Approach the Meritocratic Ideal? A Logistic Analysis," *Economic and Social Review*, 16, 115-140.
- (1993), "Maximally Maintained Inequality: Expansion, Reform and Opportunity in Irish Education, 1921-1975," *Sociology of Education*, 66, 41-62.
- Raftery, A. E., Madigan, D. M., and Hoeting (1993), "Model Selection and Accounting for Model Uncertainty in Linear Regression Models," Technical Report 262, University of Washington, Dept. of Statistics.
- Raftery, A. E., Madigan, D., and Volinsky, C. T. (1995), "Accounting for Model Uncertainty in Survival Analysis Improves Predictive Performance (with discussion)," in *Bayesian Statistics 5*, eds. J. M. Bernardo et al., Oxford, U.K.: Oxford University Press.
- Raftery, A. E., and Zeh, J. E. (1993), "Estimation of Bowhead Whale, *Balaena mysticetus*, Population Size" (with discussion), in *Bayesian Statistics in Science and Technology: Case Studies*, eds. C. Gatsonis, J. S. Hodges, R. E. Kass, and N. D. Singpurwalla, New York: Springer-Verlag, pp. 163-240.
- Rissanen, J. (1987), "Stochastic Complexity," *Journal of the Royal Statistical Society*, Ser. B, 49, 223-239.
- Robert, C. P. (1993), "A Note on Jeffreys-Lindley Paradox," *Statistica Sinica*, 3, 601-608.
- Robert, C. P., and Caron, N. (1991), "Noninformative Bayesian Testing and Neutral Bayes Factors," Technical Report 148, Laboratoire de Statistique Théorique et Appliquée, Université de Paris VI.
- Rosenkranz, S. (1992), "The Bayes Factor for Model Evaluation in a Hierarchical Poisson Model for Area Counts," Ph.D. dissertation, University of Washington, 1992, Dept. of Biostatistics.
- Rossi, P. E. (1985), "Comparison of Alternative Functional Forms in Production," *Journal of Econometrics*, 30, 345-361.
- (1988), "Comparison of Dynamic Factor Demand Models: Dynamic Econometric Modeling," in *Proceedings of the Third International Symposium in Economic Theory & Econometrics*, eds. W. A. Barnett, E. R. Berndt, and H. White, Cambridge, U.K.: Cambridge University Press, pp. 357-376.
- Rubin, D. B. (1984), "Bayesianly Justifiable and Relevant Frequency Calculations for the Applied Statistician," *The Annals of Statistics*, 12, 1151-1172.
- Schotman, P., and van Dijk, H. (1991), "A Bayesian Analysis of the Unit Root in Real Exchange Rates," *Journal of Econometrics*, 49, 195-238.
- Schwarz, G. (1978), "Estimating the Dimension of a Model," *The Annals of Statistics*, 6, 461-464.
- Shibata, R. (1976), "Selection of the Order of an Autoregressive Model by Akaike's Information Criterion," *Biometrika*, 63, 117-126.

- (1980), "Asymptotically Efficient Selection of the Order of the Model for Estimating Parameters of a Linear Process," *The Annals of Statistics*, 8, 147–164.
- (1981), "Optimal Selection of Regression Variables," *Biometrika*, 68, 45–54.
- Sklar, R., and Strauss, B. (1980), "Role of the *uvrE* Gene Product and of Inducible O⁶-Methylguanine Removal in the Induction of Mutations by *N*-methyl-*N*-nitro-*N*-nitrosoguanidine in *Escherichia coli*," *Journal of Molecular Biology*, 143, 345–363.
- Slate, E. (1994), "Parameterizations for Natural Exponential Families With Quadratic Variance Functions," *Journal of the American Statistical Association*, 89, 1471–1482.
- Smith, A. F. M. (1991), Contribution to the discussion of "Posterior Bayes Factors," *Journal of the Royal Statistical Society*, Ser. B, 53, 132–133.
- Smith, A. F. M., and Roberts, G. O. (1993), "Bayesian Computation via the Gibbs Sampler and Related Markov Chain Monte Carlo Methods," *Journal of the Royal Statistical Society*, Ser. B, 55, 3–23.
- Smith, A. F. M., and Spiegelhalter, D. J. (1980), "Bayes Factors and Choice Criteria for Linear Models," *Journal of the Royal Statistical Society*, Ser. B, 42, 213–220.
- (1981), "Bayesian Approaches to Multivariate Structure," in *Interpreting Multivariate Data*, ed. V. Barnett, Chichester, U.K.: John Wiley.
- Smith, R. L. (1989), "Extreme Value Analysis of Environmental Time Series: An Application to Trend Detection in Ground-Level Ozone" (with discussion), *Statistical Science*, 4, 367–377.
- Spiegelhalter, D. J., and Smith, A. F. M. (1982), "Bayes Factors for Linear and Log-Linear Models With Vague Prior Information," *Journal of the Royal Statistical Society*, Ser. B, 44, 377–387.
- Stewart, L. (1987), "Hierarchical Bayesian Analysis Using Monte Carlo Integration: Computing Posterior Distributions When There are Many Possible Models," *The Statistician*, 36, 211–219.
- Stone, M. (1979), "Comments on Model Selection Criteria of Akaike and Schwarz," *Journal of the Royal Statistical Society*, Ser. B, 41, 276–278.
- Tanner, M., and Wong, W. (1987), "The Calculation of Posterior Distributions by Data Augmentation" (with discussion), *Journal of the American Statistical Association*, 82, 528–550.
- Thompson, E. A., and Wijsman, E. M. (1990), "The Gibbs Sampler on Extended Pedigrees: Monte Carlo Methods for the Genetic Analysis of Complex Traits," Technical Report 193, University of Washington, Dept. of Statistics.
- Tierney, L. (1989), "Bayesian Analysis in New-S and Lisp," in *Proceedings of the Section on Statistical Computing*, American Statistical Association.
- (1990), *Lisp-Stat*, New York: John Wiley.
- Tierney, L., and Kadane, J. B. (1986), "Accurate Approximations for Posterior Moments and Marginal Densities," *Journal of the American Statistical Association*, 81, 82–86.
- Tierney, L., Kass, R. E., and Kadane, J. B. (1989), "Fully Exponential Laplace Approximations to Expectations and Variances of Nonpositive Functions," *Journal of the American Statistical Association*, 84, 710–716.
- Verdinelli, I., and Wasserman, L. (1993a), "Bayes Factors, Nuisance Parameters and Imprecise Tests," Technical Report 570, Carnegie Mellon University, Dept. of Statistics.
- (1993b), "A Note on Generalizing the Savage–Dickey Density Ratio," Technical Report 573, Carnegie Mellon University, Dept. of Statistics.
- Worsley, K. J. (1986), "Confidence Regions and Tests for a Change-Point in a Sequence of Exponential Family Random Variables," *Biometrika*, 73, 91–104.
- Zellner, A. (1971), *An Introduction to Bayesian Inference in Econometrics*, New York: John Wiley.
- (1978), "Jeffreys–Bayes Posterior Odds Ratio and the Akaike Information Criterion for Discriminating Between Models," *Economics Letters*, 1, 337–342.
- (1984), "Posterior Odds Ratios for Regression Hypotheses: General Considerations and Some Specific Results," in *Basic Issues in Econometrics*, Chicago: University of Chicago Press, pp. 275–305.
- (1987), "Comment," *Statistical Science*, 2, 339–341.
- Zellner, A., and Siow (1980), "Posterior Odds for Selected Regression Hypotheses," in *Bayesian Statistics 1*, eds. J. M. Bernardo et al., Valencia, Spain: Valencia University Press, pp. 585–603. Reply, pp. 638–643.