

-----  
DEPARTMENT OF STATISTICS  
-----

University of Wisconsin-Madison

DEPARTMENT OF STATISTICS  
UNIVERSITY OF WISCONSIN  
MADISON, WISCONSIN

DEC 21 1979

TECHNICAL REPORT NO. 573

June 1979

SAMPLING AND BAYES' INFERENCE  
IN THE ADVANCEMENT OF LEARNING

George E. P. Box

DEC 21 1979

SAMPLING AND BAYES' INFERENCE IN THE ADVANCEMENT OF LEARNING\*

George E. P. Box

June 1979

ABSTRACT

Scientific method is a process of guided learning in which accelerated acquisition of knowledge relevant to some question under investigation is achieved by a hierarchy of iterations in which induction and deduction are used in alternation.

This process employs a developing model (or series of models implicit or explicit) against which data can be viewed. Ideally at any given stage of an investigation, such a model approximates relevant aspects of the studied system and motivates the acquisition of further data as well as its analysis. By the use of a prior distribution it is possible to represent some aspects of such a model as completely known and others as more or less unknown.

Now parsimony requires that, at any given stage, the model is no more complex than is necessary to achieve a desirable degree of approximation and since each investigation is unique we cannot be sure in advance that any model we postulate will meet this goal. Therefore, at the various points in our investigation where data analysis is required, two types of inference are involved: model criticism and parameter estimation. To effect the latter, conditional on the plausibility of the model, and given the data, we can, using Bayes' Theorem, deduce posterior distributions for unknown parameters and so make inferences about them. But, before we can rely on such conditional deduction, we ought logically to check whether the model postulated accords with the data at all and, if not, consider how it should be modified. In practice, this question is usually investigated by inspecting residuals, by other informal techniques, and sometimes by making formal tests of goodness of fit. In any case model criticism, the inferential procedure whereby the need for model modification is induced, is ultimately dependent on sampling theory argument. These principles are formalized by an appropriate analysis of Bayes' formula, and implications for robust estimation are considered.

AMS (MOS) Subject Classifications - 62.02, 62A15, 62F05, 62F10

Key Words - Bayesian inference, Sampling theory inference, Estimation, Models, Predictive distribution

Work Unit Number 4 - Probability, Statistics, and Combinatorics

\* A paper read at the International Meeting on Bayesian Statistics,  
May 28-June 2, 1979 at Valencia, Spain.

Also issued as Technical Report No. 573, by the Department of Statistics,  
University of Wisconsin-Madison, Madison, Wisconsin 53706

Sponsored by the United States Army under Contract No. DAMC29-75-C-0024. & DAMC29-78-G-0166

SIGNIFICANCE AND EXPLANATION

Sampling theory inference (e.g. inference based on sampling distributions of statistics and in particular on significance tests) and Bayesian inference are usually thought of as rivals and much effort has been spent in propounding their relative merits. In this paper it is argued that both kinds of inference are needed in the scientific iteration whereby knowledge is acquired.

This iteration employs a directed alternation between induction and deduction which uses model criticism on the one hand and parameter estimation on the other. An analysis of Bayes' formula reveals model criticism as a sampling theory concept and parameter estimation as a Bayesian concept. The implications of these ideas for robust estimation are discussed.

The responsibility for the wording and views expressed in this descriptive summary lies with MRC, and not with the author of this report.

SAMPLING AND BAYES' INFERENCE  
IN THE ADVANCEMENT OF LEARNING

George E. P. Box

Today Statistics appears to be in a somewhat confused state\*. The controversy about Bayesian inference and Sampling Theory inference which some believe involves a critical choice is not resolved to most people's satisfaction. Furthermore concepts such as Data Analysis and Robust Estimation are receiving such new emphasis that some advocates of the "new Statistics" are even claiming that all else is useless and old hat.

To some extent the new and admirable emphasis on "looking at the data" is a reaction to previous extremes. On the one hand overemphasis on theory for its own sake (mathematistery) and on the other a knee-jerk approach to statistical analysis (cookbookery)†. Neither of these aberrations was healthy and some adjustment was long overdue. However I think the mistake continues to be made of assuming that different approaches to Statistics are necessarily in an adversary position. I will develop the contrary view and try to show how I believe the pieces fit together.

I start from the idea that Statistics is or should be the art and science of building scientific models which (necessarily) involve probability. Consider then how such stochastic model building should be done.

\*What is happening is related to the revolutionary change in computational speed. We need to be deterred less and less by the number of steps required in a calculation although we must correspondingly increase our concern that the human mind is also adequately involved in directing the tactics and strategy of investigation.

†See discussion of "mathematistery" and "cookbookery" in Science and Statistics, (Box 1976).

Sponsored by the United States Army under Contract No. DAMC29-75-C-0024 and DAMC29-78-G-0166.

1. The advancement of learning as an iteration between theory and practice

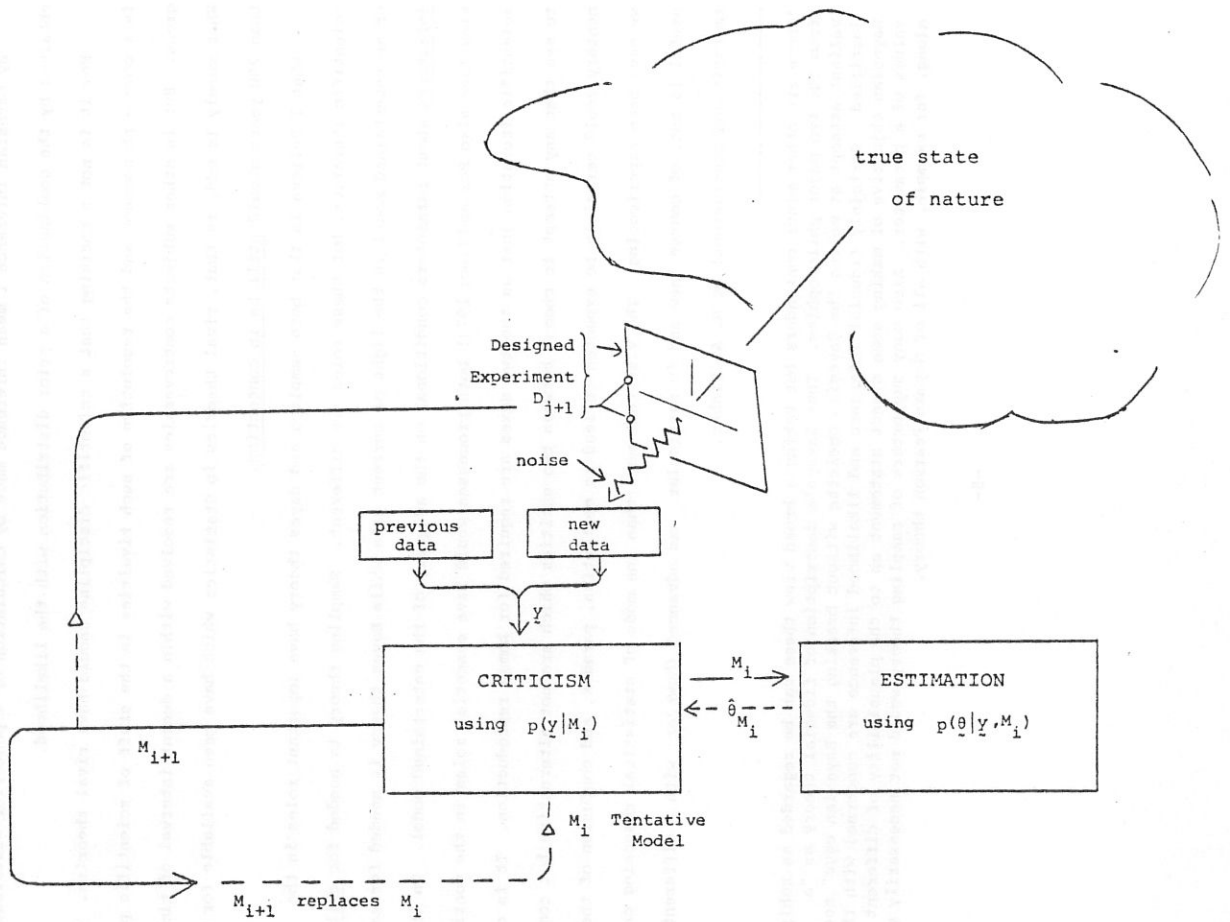
Although the matter was over the centuries debated it seems long ago to have been agreed that scientific knowledge is efficiently advanced, not by mere theoretical speculation on the one hand, nor by the mere accumulation of empirical facts on the other, but by a motivated iteration between these two activities. In this practice-theory iteration a tentative theory or model suggests a particular examination and analysis of data already existing or to be acquired\*. The results of this examination will then frequently suggest a modified model requiring further practical illumination and so on. The advancement of knowledge thus occurs as the result of an interplay between dual processes of induction and deduction which carry forward an iteration in which the model is not fixed but is continually changing. The statistician's role is to assist this process. In doing so he uses two inferential devices that I will call Criticism† and Estimation. The first can induce model modification, the second leads to estimation of unknown parameters assuming the truth of the model. For illustration, in Figure 1 at some stage of an investigation, model  $M_i$  is currently being entertained.

Criticism involves a confrontation of  $M_i$  with available data  $y$  and asks whether  $M_i$  is consonant with  $y$  and, if not, how not. It is a process of diagnostic checking. It may be done informally using plotting techniques of various kinds often involving residual quantities and more formally, with tests of goodness of fit. It may suggest that model modification to  $M_{i+1}$  is needed. In some instances it will be judged appropriate to now confront  $M_{i+1}$  with the same data, in others the nature of the modified model or necessity for independent verification may indicate the need for new data generated by a new design  $D_{j+1}$ . This will be chosen to explore shadowy regions whose illumination is currently believed to be important to progress.

Estimation. If the process outlined above leads to a verifiable model, that is one which when put to the test appears to provide an adequate approximation to reality, it may logically be used to estimate parameters conditional on its truth. However in practice this

\*I shall suppose that data is acquired from a designed experiment but the same argument would apply if data acquisition was from a sample survey or even from a visit to the library.

†The apt naming of model criticism is due to Cuthbert Daniel.



estimation process will be used not only at the termination of the model building sequence but at every stage throughout it. This is because in order to conduct criticism of the model it is often necessary to provisionally estimate parameters at intermediate stages, tentatively entertaining the model as if it were believed true.

I shall argue in this paper that while criticism must ultimately appeal to sampling theory for its justification estimation requires the use of Bayes theorem (or, for the faint-hearted, likelihood). Acceptance of this position provides justification for a specific kind of sampling theory significance tests but none for sampling theory confidence intervals.

## 2. Rival theories of inference

The distinction between inferential criticism and parameter estimation has often not been made and proponents both of sampling inference and Bayesian inference\* have long sought, mistakenly in my view, for a single comprehensive theory. By sampling theory inference I mean inference made by referring some relevant function of the data to a reference sampling distribution which would be appropriate if some specific hypothetical model  $M_0$  were true.

By Bayesian inference I mean inference made by calculation of a posterior distribution obtained by the combination of a prior distribution with the likelihood.

Now it is not surprising that a scientific discipline should have rival theories. This is a common phenomenon and the resolution of such rivalries is the stuff of scientific progress. But in other subjects controversies are resolved within a decent interval of time. What surely is odd, is that, rival theories in Statistics which have been available for more than 200 years should still be in contention.

What I believe is that both sampling and Bayes theory have important roles in the scientific iteration, but these roles are different. Sampling theory is needed for criticism of an entertained model in the light of current data while Bayes theory is needed for making inferences about parameters conditional on the adequacy of the entertained model. On this view (see also Box and Tiao; 1973) both processes would have essential roles in the continuing scientific iteration just as the two sexes are required for human reproduction. It is easy to see that any attempt to choose between two entities which are not alternative but complementary could certainly be expected to lead to contention, paradox, and confusion of the kind we have been experiencing. The view that more than one mode of statistical reasoning can be useful is not, of course, new and in particular was advanced (however with a different emphasis and conclusions) by R. A. Fisher.

\*There are other minor contenders but taking a broad view these can be regarded as schisms from the two major philosophies. Thus Savage's description of fiducial theory as "a valiant attempt at making the Bayesian omelette without breaking the Bayesian eggs" seems justified. Certainly fiducial inference and likelihood inference are concerned with the Bayesian objective of making some direct statement as to the plausibility of different values of a parameter. Also many supporters of sampling theory would not necessarily go along, for example, with all of Neyman-Pearson theory.

## 3. Some remarks on Sampling and Bayes inference

The essence of what I mean by "sampling theory inference" is exemplified by the Shewhart quality control chart. The set of limit lines for the sample mean for example indicates for this function of the data, a reference distribution appropriate for the model  $M_0$  (for the process in control). A single outlying point is surprising because it is associated with unusually low probability density. It thus raises the possibility that  $M_0$  is inappropriate and that an alternative model  $M_1$  might be needed to explain the inadequacy. In the words of Shewhart, the process is out of control in a manner which we may be able to attribute to an assignable cause. A number of different functions of the data may be considered in checking the appropriateness of  $M_0$  and their nature depends on the type of departures from  $M_0$  that are in mind. Thus quality control charts are often kept of both the sample mean and the sample range to indicate departures from  $M_0$  in both level and spread and other functions such as run length of positive deviations might also be considered. Finally patterns which were not foreseen may possibly turn up, invite consideration, and induce possible explanations to be subsequently tested.

### Prior probabilities in Bayesian and Sampling inference

In the past the need for prior probabilities has often not been thought of as a necessity for all scientific inference but rather as a feature peculiar to Bayesian inference. Indeed it is often regarded by non-Bayesians as the major point of weakness of Bayes theory and has, therefore, been a focus for attack and sometimes for derision. By contrast a Bayesian proponent might argue (a) that any theory of estimation worthy of the name should make it possible, given a model, to say after data had come to hand what was believed about the values of its parameters and (b) that what was believed after the data was available must surely depend on what was believed before it was available (c) that this would include the possibility of sometimes using non-informative prior distributions either to simulate the actual state of relative ignorance of the investigator or to represent the impact of the data on a hypothetical unbiased observer (or juror). He might argue further that the difficulties and paradoxes that have embarrassed advocates of sampling theory as it

has been practiced and their inability to fix up the theory convincingly have come from its past inadequate capability to include prior information.

Sampling theory is of course not free from assumptions of prior knowledge. Instead it is as if only two states of mind have been allowed--complete certainty or complete uncertainty. Whereas in the sampling theory context a parameter had to be treated either as exactly known or as completely unknown, in the Bayesian context a prior could be chosen to approach either of these extremes or any intermediate state.

In this connection it is important to remember that every simple model can be thought of as embedded in a more complex one. For example an outright assumption of normality can be modelled by a suitable parametric family of distributions indexed by a parameter  $\theta$ , which has a sharp prior at the normal value. Independence of errors, so frequently assumed, can similarly be represented by a sharp prior operating on a broader model allowing appropriate dependence. Seen in this way, it appears that, when assumptions of normality and independence are made in sampling theory, it is not that no prior knowledge is used, but rather that implausibly precise prior knowledge is implied.

#### 4. The model is the prior

Such considerations lead me to believe that it is impossible to logically distinguish between the model and the prior distribution. In a real sense the model is the prior. A model is a probability statement of all the assumptions currently to be tentatively entertained a priori. These probability statements can express certainty or various degrees of uncertainty.

Of course models are approximations (good ones are artfully chosen approximations which work well in practice). But there is good reason to believe that the "all or none" prior assumptions implied in the traditional use of a sampling theory are frequently too crude even as an approximation. Indeed many of the difficulties of sampling theory which have come to light in recent years may be traced to the primitive means it has available for incorporating prior knowledge and the crippling effect of allowing only probability statement of a certain kind to be included in the model. One illustration of how implied prior knowledge which is implausibly imprecise can lead to trouble in sampling theory in the famous discovery by Stein (1956) of the inadmissibility of normal multivariate mean, and the improved nonlinear shrinkage estimators which give smaller mean square error.

It is however easy to miss the lesson which is to be learned from such examples. To be specific, consider the usual one-way analysis of variance set-up. Here a locally uniform prior distribution for the set of group means  $\mu' = (\mu_1, \mu_2, \dots, \mu_j, \dots, \mu_n)$  which would exactly justify the sample averages as estimators makes little sense (see, for example, Box and Tiao (1968), Lindley and Smith (1972)). By contrast the prior assumption which justifies the shrinkage estimator is that the  $\mu_j$  are random samples from some normal super population having unknown mean and variance. This corresponds to the usual "model II" sampling theory assumption and in appropriate circumstances could be eminently reasonable. It is crucial to notice, however, that there are many circumstances in which this latter assumption would not be sensible either because, although prior knowledge about  $\mu_1, \mu_2, \dots, \mu_n$  existed, it was of quite a different character. For example, if the  $\mu$ 's were daily batch yields from some production process, it would usually be much more sensible to



Postulate that the allowed some time series model such as a stationary autoregressive process (Tiao and Ali (1971)). The estimators then derived from Bayesian means are not Stein's shrinkage estimators, which would appropriately introduce sample information about  $\sigma^2$ , but alternative estimators allowing incorporation of relevant sample information about the autocorrelation of the batch means.

Some sampling theorists concede that Bayes theorem may be used as a kind of conjuring trick to produce efficient estimators which are then used in a sampling theory context. In this exercise they regard the prior distribution as a convenient prop which is never taken seriously and is quickly discarded. I think the example quoted above is one of many which shows that this idea has no rational status. For it illustrates that there is not one set of "shrinkage estimators", but an infinity of such sets depending (very naturally) on the model (that is the prior) which is appropriate to describe the particular scientific situation under study.

The strength of the explicit statement of prior assumptions is that in the iterative model building process, they make manifest at every stage exactly what assumptions are tentatively entertained and so allow them to be criticized. Some of the nervousness experienced by non-Bayesians confronted with the idea of a prior distribution has perhaps arisen because the iterative nature of scientific process and consequent tentative transitory character of models and all their assumptions, has not been generally understood.

Many of us were taught to think unrealistically in terms of "one shot" procedures.

The sequence: frame hypothesis - collect data - test hypothesis/make decision; of course, fails to describe the usual context in which Statistics is applied.

Critics have therefore feared gross mistakes arising from adamant prior prejudice which ignored "what the data were trying to say." In the iterative context of real scientific enquiry however gross mistakes about the prior or any other aspect of the model will usually be corrected at the criticism phase of the next iteration.

##### 5. Two complementary factors from Bayes formula

If we accept the prior probability distribution of parameters  $\theta$  as an essential part of a model then all aspects of the model, hypothesized at some particular stage of an investigation, are contained in the joint density obtained by combining the likelihood and the prior

$$P(Y, \theta | M) = P(Y | \theta, M) \cdot P(\theta | M) \quad (5.1)$$

where  $|M$  is understood to indicate conditionality on some aspect of the model and  $Y$  is a data vector.

This joint distribution which is a comprehensive statement of the model can also be factored as

$$P(Y, \theta | M) = P(\theta | Y, M) P(Y | M) \quad (5.2)$$

and can be computed before any data become available. In particular the second factor on the right

$$P(Y | M) = \int P(Y | \theta, M) P(\theta | M) d\theta, \quad (5.3)$$

which is the predictive distribution, may be so calculated. It is the distribution of the totality of all possible samples that could occur if the model  $M$  were true.

When an actual data vector  $Y_d$  becomes available

$$P(Y_d, \theta | M) = P(\theta | Y_d, M) P(Y_d | M). \quad (5.4)$$

The first factor on the right is then Bayes' posterior distribution of  $\theta$  given  $Y_d$

$$P(\theta | Y_d, M) = k P(Y_d | \theta, M) P(\theta | M) \quad (5.5)$$

and the second factor

$$P(Y_d | M) = \int P(Y_d | \theta, M) P(\theta | M) d\theta = k^{-1} \quad (5.6)$$

is the predictive density associated with the data set  $Y_d$  actually obtained. Figure 2 illustrates for a single parameter  $\theta$  and a sample  $Y_d$  of  $n = 2$  observations.

If the model is to be believed, then the posterior distribution  $P(\theta | Y_d, M)$  allows all relevant estimation inferences to be made about  $\theta$ . However even if the model were totally incorrect, this could not be shown by any abnormality in this factor which is conditional on both data and model specification. However plausibility or otherwise of obtaining such a sample if the model were appropriate may be assessed by reference of the density  $P(Y_d | M)$  to

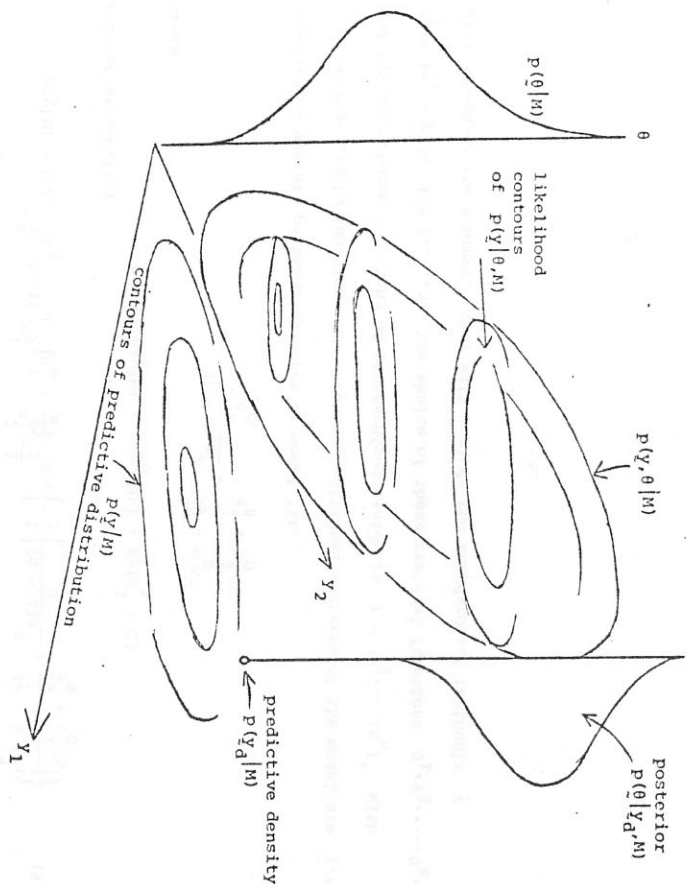


Figure 2. Showing for a single parameter  $\theta$  and sample  $\bar{y}_d$  of two observations; the prior distribution, likelihood contours, the posterior distribution and contours of the predictive distribution.

the predictive reference distribution  $p(y|M)$ . An unusually small value of  $p(y_d|M)$  as measured by  $1/p(y_d|M) < p(y_d|M)$  casts doubt on the appropriateness of the model  $M$ . Now  $p(y|M)$  is an  $n$ -dimensional distribution and it will usually be true that if the model is inadequate it is most likely to be deficient in certain directions associated with unusual values of certain specific functions  $g_1(y)$  of the data. Examples of such functions are sample averages, variances, moment coefficients, coefficients of serial correlation, and other measures of standardized deviations from a norm. In every case the appropriate reference distribution to which the realized statistic  $g_1(y_d)$  should be referred is the distribution  $p(g_1(y|M))$ , when the model  $M$  is true, derived by appropriate integration of  $p(y|M)$ .

In practice, criticism or diagnostic checking of the model is often conducted by visual inspection of residual displays and other more sophisticated plots. But such a process, although it is informal, still, it seems to me, falls within the logical framework described above. The statistician is looking for "features" in the data which would be surprising or "unusual" if the model  $M$  were true. Such a feature can be described by a function  $g(y_d)$  and its unusualness, if formalized, would have to be measured by reference to  $p(g(y_d)|M)$ .

In addition to possible discrepancies to which we have been alerted by experience, other features may appear pointing to inadequacies of a kind not previously suspected. This possibility has sometimes proved perplexing for statisticians, for while on the one hand the truly unexpected could point the way to precious new knowledge, on the other, associated probabilities will be indeterminate because of the uncountable character of other features that might also have been regarded as surprising. I think the calculation which ignores this difficulty of indeterminate selection should still be made, for while it might lead to the too frequent pursuit of nonexistent assignable causes, the iterative process will quickly terminate this chase and carrying out the exercise will at least eliminate phenomena, which at first sight look surprising but really are not. For example, Feller (1968) shows that for a random group of 30 people, the probability that at least two have coincident birthdays is over 70%, this tells us we need look no further for an explanation when we are surprised to find two such people at a party.



Example: Unknown mean  $\theta$ , variance  $\sigma_0^2$  assumed known.

Consider a sample of  $n$  observations drawn randomly from a normal distribution with unknown mean  $\theta$  and known variance  $\sigma_0^2$ . We express uncertainty about the mean by supposing that a priori  $\theta$  is distributed normally about  $\bar{\theta}_0$  with variance  $\sigma_\theta^2$ .

Thus

$$p(\bar{y}|\theta, N) = (2\pi)^{-\frac{n}{2}} \sigma_0^{-n} \exp\left\{-\frac{1}{2} \frac{\sum_{i=1}^n (y_i - \theta)^2}{\sigma_0^2}\right\} \quad (5.7)$$

$$p(\theta|N) = (2\pi)^{-\frac{1}{2}} \sigma_\theta^{-1} \exp\left\{-\frac{1}{2} \frac{(\theta - \bar{\theta}_0)^2}{\sigma_\theta^2}\right\} \quad (5.8)$$

The posterior distribution from which  $\theta$  may be estimated conditional on the adequacy of the model is then

$$p(\theta|\bar{y}, N) = (2\pi)^{-\frac{1}{2}} \left\{ \frac{1}{\sigma_0^2} + \frac{n}{\sigma_\theta^2} \right\}^{-\frac{1}{2}} \exp\left\{-\frac{1}{2} \left[ \frac{1}{\sigma_0^2} + \frac{n}{\sigma_\theta^2} \right] (\theta - \bar{\theta})^2\right\} \quad (5.9)$$

where  $\bar{\theta} = \left( \frac{1}{\sigma_0^2} + \frac{n}{\sigma_\theta^2} \right)^{-1} \left( \frac{1}{\sigma_0^2} \bar{\theta}_0 + \frac{n}{\sigma_\theta^2} \bar{y} \right)$

The predictive distribution which can act as reference distribution for the observed data vector  $\bar{y}_d$ , thus allowing criticism of the model, is

$$p(\bar{y}|N) = (2\pi)^{-\frac{n}{2}} \sigma_0^{-(n-1)} \frac{1}{n} \frac{\sigma_0^2}{\sigma_\theta^2} \exp\left\{-\frac{1}{2} \left[ \frac{(n-1)s^2}{\sigma_0^2} + \frac{(\bar{y} - \bar{\theta}_0)^2}{\sigma_\theta^2/n} \right]\right\} \quad (5.10)$$

And the probability

$$P = \Pr(p(\bar{y}|N) < p(\bar{y}_d|N)) = \Pr\{X_n^2 > C\}$$

where

$$C = \frac{(n-1)s^2}{\sigma_0^2} + \frac{(\bar{y} - \bar{\theta}_0)^2}{\sigma_\theta^2/n}$$

supplies an overall portmanteau check on model fit.

Obvious sample functions for checking individual features of the model are  $\bar{y}, s^2$  and suitably chosen functions of standardized residuals  $\bar{r} = (r_1, \dots, r_n)'$  with

$r_i = (y_i - \bar{y})/s$   $i = 1, \dots, n$ . The choice of these residual functions  $g_1, g_2, \dots, g_k$  will depend on the context. They will include the standardized residuals  $\bar{r}$

themselves, but might also address the need to apply checks for "bad values", skewness, kurtosis and serial correlation, for example. The standardized residuals in the form defined above are constrained by the identities  $\sum r_i = 0$ ,  $\sum r_i^2 = n-1$  and can be more conveniently parameterized in terms of  $n-2$  independently distributed functions obtained as follows:

Make an orthogonal transformation from  $\bar{y}$  to  $\bar{y} = (y_1, y_2, \dots, y_n)'$  with  $y_n = \sqrt{n}\bar{y}$  and then transform to  $\bar{y}, s^2$  and  $\bar{u}$  where  $\bar{u}$  is a vector of  $n-2$  residual quantities  $\bar{u} = (u_1, u_2, \dots, u_{n-2})'$  such that

$$u_j = y_{j+1} / \left\{ \sum_{i=1}^j y_i^2 / j \right\}^{\frac{1}{2}}.$$

The Jacobian of the transformation from  $\bar{y}$  to  $\bar{y}, s^2, \bar{u}$  is proportional to

$$\frac{n-1}{s^2} \prod_{j=1}^{n-2} (1 + u_j^2/j)^{-\frac{1}{2}(j+1)}$$

. After transformation therefore the predictive distribution contains  $n$  elements all of which are distributed independently and becomes

$$p(\bar{y}, s^2, \bar{u}|N) = p(\bar{y}|N) p(s^2|N) p(\bar{u}|N) \quad (5.11)$$

where

$$p(\bar{y}|N) = (2\pi)^{-\frac{n-1}{2}} (\sigma_0^2 + \sigma_\theta^2/n)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2} (\bar{y} - \bar{\theta}_0)^2 / (\sigma_0^2 + \sigma_\theta^2/n)\right\} \quad (5.12)$$

$$p(s^2|N) = \left\{ \frac{1}{2} (n-1) \right\}^{\frac{1}{2}} \Gamma^{-\frac{1}{2}} \left\{ \frac{1}{2} (n-1) \right\} (\sigma_0^2)^{-\frac{1}{2}} \frac{n-1}{s^2} \exp\left\{-\frac{1}{2} (n-1) s^2 / \sigma_0^2\right\} \quad (5.13)$$

$$p(\bar{u}|N) = \Gamma^{-\frac{1}{2}} (n-1) \frac{(n-1)}{2} \prod_{j=1}^{\frac{1}{2}(n-1)} \left\{ 1 + \frac{u_j^2}{j} \right\}^{-\frac{1}{2}(j+1)} \quad (5.14)$$

The standardized residual quantities of interest  $g_1, g_2, \dots, g_k$  can then be expressed equally as functions  $f_1(u), f_2(u), \dots, f_k(u)$  of the  $u$ 's. So that, in particular, unusual features of  $\bar{y}, s^2$  and  $g_1, \dots, g_k$  given the model could be assessed by

computing

$$(i) \Pr(p(\bar{y}|N) < p(\bar{y}_d|N))$$

$$(ii) \Pr(p(s^2|N) < p(s_d^2|N))$$

$$(iii) \Pr(p(g_j|N) < p(g_{jd}|N)) \quad j = 1, 2, \dots, k.$$

these are the (two tail area) probabilities associated with reference of

- (i)  $(\bar{y}_d - \theta_0)/(s_0^2/n)^{1/2}$  to the Normal table
- (ii)  $(n-1)s_d^2/s_0^2$  to the  $\chi^2$  table
- (iii)  $g_d$  to the reference distribution obtained by appropriate integration of the distribution  $p(\bar{y}|M)$ .

They yield checks on the adequacy of the model which we denote by  $c(\bar{y}), c(s^2), c(g)$ .

For example suppose the yield of a batch process was under study and that a sample  $\bar{y}$  was available of  $n$  observations all from a single batch having unknown mean  $\theta$ . Suppose at this stage of the investigation that the tentative model assumed that, because of process variation, batch means varied Normally and independently about some value  $\theta_0$  with variance  $\sigma_0^2$  and, because of testing variation, the  $i$ th observation  $y_i$  varied about  $\theta$  normally and independently with variance  $\sigma_0^2$ . Then the model would be that discussed above and, if this model could be believed, the batch mean  $\theta$  would be estimated by the posterior distribution  $N(\bar{y}, (I_{\bar{y}} + I_y)^{-1})$  where  $I_{\bar{y}} = n\sigma_0^{-2}$ ,  $I_y = \sigma_0^{-2}$ . And, if we write  $w = I_{\bar{y}}/(I_{\bar{y}} + I_y)$  for the proportion of the information coming from the sample, then  $\bar{\theta} = w\bar{y} + (1-w)\theta_0$ .

Before drawing such a conclusion however a prudent statistician would question the

- (i) In particular applying the checks  $c(\bar{y}), c(s^2), c(g)$ ,
- (ii) an unusually small value of  $p(\bar{y}|M)$  could call into question the choice of some or all of  $\theta_0, \sigma_0^2$  and  $\sigma_0^2$ .
- (iii) an unusually small value of  $p(s^2|M)$  could call into question the choice of  $\sigma_0^2$ .
- (iii) an unusually small value of  $p(g_j|M)$  could suggest departures from the assumed distributional form  $p(y|\theta, M)$  produced by serial correlation, bad data values, non-normality, etc.

Only after the investigator had found that the evidence offered by the data did not invalidate the model should he proceed to make the conditional deductive inference supplied by Bayes theorem.

## 6. Some Implications

Consider the problem of making inferences about  $\theta$  in the previous example. If we assume the model true then we can estimate  $\theta$  from a normal posterior distribution with mean  $\bar{\theta} = w\bar{y} + (1-w)\theta_0$  and variance  $(I_{\bar{y}} + I_y)^{-1}$  where  $w = I_{\bar{y}}/(I_{\bar{y}} + I_y)$  is the fraction of the information coming from the sample. First however we require to check the model using the predictive distribution. In particular the check  $c(\bar{y})$  requires a reference of  $(\bar{y}_d - \theta_0)/(s_0^2/n)^{1/2}$  to the normal table. Significance test. Suppose  $\sigma_0^2$  is assumed small compared with  $\sigma_0^2/n$ , then  $w$ , the relative amount of information, supplied by the data is small and  $1-w$  is close to unity. Then, if this model can be relied upon, the posterior distribution is essentially the same as the prior and is sharply centered at  $\theta_0$ . (A practical context is one where the statistician is told that process variation is negligible compared with testing variation and the process mean is known to be  $\theta_0$ .) If this model is assumed, then information from available data  $\bar{y}$  can add very little to what is known already. However, it can deny the relevance of this model. In particular  $c(\bar{y})$  involves the reference of  $(\bar{y} - \theta_0)/(s_0^2/n)^{1/2}$  to normal tables; the failure of this check means that the model is discredited and therefore the operation that leads to a sharp posterior distribution centered at  $\theta_0$  may not logically be undertaken.

The above most satisfactorily explains to me the rationale of a significance test.

- (i) The tentative model (null hypothesis) implies that  $\theta = \theta_0$ .
- (ii) A check on this aspect of the model is provided by reference of  $(\bar{y} - \theta_0)/(s_0^2/n)^{1/2}$  to the Normal Table.
- (iii) If the tail area probability is not small we do not question the model. The application of Bayes theorem then produces a posterior distribution which is a delta function at  $\theta_0$ . We have "no reason to question the null hypothesis".
- (iv) If the tail area probability is small we conclude that the model which postulates that  $\theta = \theta_0$  is discredited by the data and that some other model is appropriate. The "null hypothesis is rejected."

(v) Notice too that although the failure of this check would not immediately prescribe the use of Bayes theorem, the failure of other checks (and of  $c(s^2)$  in particular) would also indicate the necessity of model modification before proceeding further.

A difficulty that this removes for me is that, as usually formulated, significance tests seem to provide no basis for belief. On the above argument, if we accept the model, we believe a priori that  $\theta$  is close to  $\theta_0$ . We must therefore believe that  $\theta = \theta_0$  very nearly a posteriori. The availability of data provides however an opportunity to assess the concordance of data and model.

The significance test itself provides a means only of discrediting the model. Our belief in the proposition  $\theta = \theta_0$  comes from an application of Bayes theorem for a model which there is no reason to question (as a reasonable approximation to truth).

In particular this underscores the illogicality of testing a null hypothesis which is not credible to begin with. Thus the Durbin-Watson test for serial correlation, for which the null hypothesis is that errors are distributed independently, is frequently misapplied to test serial data which a priori can be expected to be autocorrelated.

#### Precise measurement and improper priors

Suppose now that  $\sigma_\theta^2$  was very large compared with  $\sigma_0^2/n$ . The predictive check  $c(\bar{y})$  now approaches  $(\bar{y}_d - \theta_0)/\sigma_\theta$  implying that for sets of data having widely different sample averages the model would not be called into question. The situation where such a non-informative prior distribution was relevant was referred to by L. J. Savage as that where the theory of precise measurement applied. The invocation of this principle might, at first, seem a license to use Bayes theorem without any restraining checks of the model. But this idea makes no sense either from an applied or a theoretical point of view.

The practical situation is that the sample information coming from  $\bar{y}$  must be evaluated in a context where there is relatively very little prior information about the value of  $\theta$ .

Here computational convenience and logic must of course be carefully distinguished. Replacing "relatively very little" by zero can be justified computationally in those circumstances where to do so provides a good numerical approximation but not otherwise. However in either case zero remains infinitely smaller than any small quantity. In this example, substitution of an improper uniform prior will produce a normal posterior distribution having mean  $\bar{y}$  and variance  $\sigma_0^2/n$ , also obtained as the limit when in our model, the fraction of information  $w$  supplied by the data tends to unity. But not only that, the specification of the prior for  $\theta$  as  $N(\theta_0, \sigma_\theta^2)$  is obviously overly specific, and the improper prior could provide an appropriate limit for disperse priors which were widely different in structure and/or much less specific.

All statistical results, in so far as they relate to reality, are approximations. Those obtained from improper priors do in many important examples provide excellent approximations. I hasten to add of course that limiting processes can be tricky and theoretical statisticians are right to worry about them.

Notice however that the situation is different for the predictive check. To say that  $w$  is close to unity is only to say that  $\sigma_\theta^2$  will dominate the denominator in  $(\bar{y} - \theta_0)/(\sigma_\theta^2 + \frac{\sigma_0^2}{n})^{\frac{1}{2}}$ . But to say that it is equal to unity implies that  $\sigma_\theta^2$  is infinite and the check cannot be made, which implies that there are absolutely no values of  $\bar{y}$  which could discredit the model - a situation which I cannot imagine as practically possible.

Consider for example, a physical chemist who runs experiments to determine the activation energy  $\theta$  for a particular chemical reaction about which little is known. It would usually be true that his initial uncertainty about  $\theta$  would be large compared with the anticipated standard deviation  $\sigma/\sqrt{n}$  of the experimental procedure, the theory of precise measurement would apply therefore and the limiting result obtained from the usual improper prior would supply a good approximation. Nevertheless the chemist may know that activation energies for compounds of the kind being tested are usually measured in tens of kilo calories per gram mole. If the statistician, who has perhaps misplaced a decimal point, presents him with an estimate of

say 0.1 kilo calories per gram mole he will rightly reject it. In doing so he will be informally conducting a check formalized by  $c(\bar{y})$ . In practice then checks such as  $c(\bar{y})$  can never really be dispensed with. The non-informative prior used in practice must to make practical sense always be proper, but nevertheless the appropriate posterior distribution can, in suitable circumstances, be numerically approximated by the device of substituting an improper prior. I labour this point because although it has been made earlier (see for example Box and Tiao 1973, p. 28) critics seem to have misunderstood earlier discussions. Explicit consideration of predictive checks makes the situation even clearer.

#### Choosing the diagnostic checks

Frequently the checking functions  $q(\bar{y})$  which are to be used formally or informally for checking various features of a model  $M$  are chosen on an ad hoc basis.

One formal basis for selection of such functions follows essentially the route explored by Neyman and Pearson. Suppose a basic model  $M_0$  is given and an alternative model  $M_1$  represents some discrepancy from  $M_0$  which is of interest. Then a function of the data suitable for detecting such discrepancies may be obtained from the ratio<sup>†</sup>

$$p(y_d|M_0)/p(y_d|M_1)$$

#### Parsimony: Diagnostic checks versus Robustification

A question which confronts\* the statistician at every stage of an investigation is "How complex a model should I use?" The possibilities for model elaboration are of course limitless. For instance a commonly used model assumes errors to be independently, identically and Normally distributed (IIN). It is easy to imagine a sequence of fall-back models which might begin like this

$$M_0 + M_1 + M_2 + M_3 + \dots$$

$$IIN \quad IIN \quad IIN \quad IIN$$

<sup>†</sup>Model criticism cannot logically be conducted by the study of the magnitude of such ratios however, for even if this ratio were very high the predictive check could still show the favored model to be highly implausible.

\*An apparently different question is "Should I use a robust procedure?", but I will argue that this is subsumed by the broader question.

At each stage of elaboration there are many forms the modified model could take and most require additional parameter values either to be given from prior knowledge or to be estimated from the data. Obviously compromise is necessary, for, on the one hand, simpler models can allow better scientific understanding and better estimation, while, on the other hand, more complex ones can, but need not, be closer to the truth. A conclusion is that, realistically, model building is iterative, so that mistakes can be rectified.

This fact of necessary compromise raises the dilemma of where should the compromise be made, that is to say, of what should be left out and what be included. In particular suppose some deviation from an "ideal" model  $M_0$  can be parametrized by a discrepancy parameter  $\beta$  or a vector of such parameters.

For illustration  $M_0$  might be the usual normal model and  $\beta$  could measure

- (i) possible serial correlation of errors  $(e_{t-1}, e_t, e_{t+1}, \dots)$ ; for instance, the serial correlation might be generated by a first order autoregressive process  $e_t = \beta e_{t-1} + a_t$  where  $a_t$  was a source of discrete white noise.

- (ii) possible deviation from error normality; for example\* according to  $p(e|\sigma, \beta) = \text{const } \sigma^{-1} \exp[-\frac{1}{2}(e^2/\sigma^2)1/(1+\beta)]$

- (iii) need for parametric transformation; for example the normal linear model would be valid not for  $y$  but for  $y^\beta$ .

- (iv) need to allow for bad values; for example with probability  $\beta$  (close to unity) the error variance was  $\sigma^2$ , with probability  $1 - \beta$  it was  $k^2\sigma^2$ .

In each case there are two ways to handle the possible model discrepancy, depending on whether the parameter  $\beta$  is omitted from or included in the model. We call these

#### diagnostic checking and robustification.

Diagnostic checking. If the discrepancy parameter is omitted from the model then an appropriate diagnostic check can be made. Formally this would be done by referring

\*Here and elsewhere other functional forms might be found more appropriate. These examples are intended only to illustrate essential principles; not, of course, to be comprehensive.

some suitable function  $g(y)$  of the data to a reference distribution derived for the predictive distribution  $p(y|M_0)$ .

Robustification. If the discrepancy parameter is included then robust estimation\* of  $\hat{\theta}$  is provided by the posterior distribution

$$p(\hat{\theta}|y) = \int p(\hat{\theta}|\beta, y)p(\beta|\hat{y})d\beta \quad (6.1)$$

If we write

$$p_u(\beta|\hat{y}) = p(\beta|\hat{y})/p(\beta) \quad (6.2)$$

$$p(\hat{\theta}|\hat{y}) = \int p(\hat{\theta}|\beta, \hat{y})p_u(\beta|\hat{y})p(\beta)d\beta \quad (6.3)$$

In this last expression

- (i)  $p(\beta)$  can be chosen to represent approximately the probability of occurrence of different values of  $\beta$  in the real world
- (ii) the function  $p_u(\beta|\hat{y})$  is a pseudo-likelihood which reflects information about  $\beta$  supplied by the data
- (iii) considered as a function of  $\beta$ ,  $p(\hat{\theta}|\beta, \hat{y})$  reflects the sensitivity of estimation to the choice of the discrepancy parameter.

The omission of the parameter  $\beta$  is equivalent to setting it equal to the value  $\beta_0$  which it takes in the ideal model  $M_0$ . Table 1 shows some examples of diagnostic checks and corresponding robust estimation methods. A fuller discussion is given elsewhere (Box 1979).

Discussion. There may be Bayesians who would deny the need for diagnostic checks based on sampling theory. They may feel that "they can do it all with Bayes". I do not believe this position can be sustained because it implies either

- (i) that they know what the model is in advance or
- (ii) that they are prepared to make the model so comprehensive that nothing could possibly be overlooked.

\*Numerous authors (Huber, Tukey, Andrews, Hampel, etc.) have proposed ad hoc methods of robust estimation relying on the empirical modification of classical estimation procedures. It seems more logical to me to modify the model which is presumably at fault rather than the method of estimation which is not. Furthermore this has the advantage of clearly revealing the assumptions which are being made.

Example	$M_0$		$M_1$
	Make Inferences Using $p(\theta y_D, \beta = \beta_0)$	Check Using $g(y_D)$	Make Robust Inference <sup>†</sup> Using $p(\theta y_D) = \int p(\theta \beta, y_D)p(\beta y_D)d\beta$
Serial Correlation	Normal Linear Model	e.g. Durbin-Watson (1950, 1951) check	e.g. Zellner and Tiao (1964)
Kurtic Error Distribution	Normal Linear Model	e.g. Anscombe and Tukey (1963) checks for kurtosis	e.g. Box and Tiao (1962)
Transformation	Normal Linear Model	e.g. Tukey's (1949) one degree of freedom for transformation	e.g. Box and Cox (1964)
Bad Observations	Normal Linear Model	Tests for outliers e.g. Grubbs (1950), Dixon (1950), Ferguson (1961), David, Hartley and Pearson (1954)	e.g. Box and Tiao (1968)*

TABLE 1 SOME EXAMPLES OF DIAGNOSTIC CHECKS AND ROBUSTIFICATION

\*It is of course only in relation to this one problem that robust estimation is usually considered and that usually from the empirical non-Bayesian approach of Huber, Tukey, Andrews, etc.

<sup>†</sup>The robust model would, of course, also be subjected to checks.



both positions are grandiose and unrealistic and the second if attempted could lead to unacceptably complicated models which would impede scientific progress.

In this connection it must be realized that looking at residuals is essentially a sampling theory procedure and is an acknowledgement of the often happy fact that an experiment might reveal more than was bargained for. To put it another way, every Bayesian statement is conditional and somewhere there has to be an anchor.

An acceptance of my theme implies of course that what is tentatively included in a model is a matter of judgement.\* However we can still look for guidelines for model building on what to tentatively include (robustify for) and what to tentatively omit (and later check for).

Obviously the need for special features in the model depends on the context, e.g.:

(a) serial data (in particular most economic and business data) cannot reasonably be expected to be represented by a model with uncorrelated errors, autocorrelation is virtually certain (temporary changes in mean and variance are also very likely in serial data), (b) data for which  $y_{\max}/y_{\min}$  is large is likely to need transformation before any simple model could apply, (c) most experimental data are liable to occasional bad values. Elaborations which are primary candidates for robustification (inclusion in the model) reflect features which might easily elude diagnostic checks and could then invalidate subsequent analysis.

Although the ad hoc robustifiers seem to have given all their attention to possible non-normality of (assumed independent) observations, an even greater source of serious trouble is autocorrelation in serial data. See for example Coen, Gomme and Kendall (1968), Box and Newbold (1971), Palleen (1977), Box and Jenkins (1970).

\*This idea that a statistician has to use scientific judgement is not a universally popular one. The objectivity of statistics like that of science does not of course mean that all statisticians (or scientists) even though capable of using the same set of tools will do equally well when using them. Just as there are good lawyers and bad lawyers, there are good statisticians and poor ones.

Another Example:  $\theta$  known,  $\sigma^2$  unknown

Suppose now we have a known mean  $\theta$  but unknown variance  $\sigma^2$ . Also suppose we express uncertainty about the variance by assuming a priori that  $\sigma^2$  is distributed about  $s_0^2$  in a scaled  $\chi^2$  distribution having  $\nu_0$  degrees of freedom. This is equivalent to assuming that a supposedly relevant estimate  $s_0^2$  of  $\sigma^2$  having  $\nu_0$  degrees of freedom is available from past data and has been assessed against a non-informative reference prior (i.e. prior to the first sample the distribution of  $\log \sigma$  was flat in the neighborhood of the likelihood). Then for a prospective sample of  $n = \nu + 1$  observations

$$\left\{ \begin{array}{l} p(y|\sigma^2, n)(\sigma^2)^{-\frac{n}{2}} \exp \left[ -\frac{1}{2} \frac{vs^2 + n(\bar{y} - \theta)^2}{\sigma^2} \right] \\ p(\sigma^2 | n)(\sigma^2)^{-\left[ \frac{\nu_0}{2} + 1 \right]} \left( \frac{\nu_0}{s_0^2} \right)^{\frac{\nu_0}{2}} \exp \left[ -\frac{1}{2} \frac{\nu_0 s_0^2}{\sigma^2} \right] \end{array} \right. \quad (A.1)$$

The complete prospective statement about the model is thus

$$p(y, \sigma^2 | n)(\sigma^2)^{-\frac{\nu_0}{2} - \left[ \frac{n + \nu_0}{2} + 1 \right]} \exp \left[ -\frac{1}{2} \frac{(\nu_0 + n)\bar{\sigma}^2}{\sigma^2} \right] \quad (A.2)$$

where  $\bar{\sigma}^2 = (n\bar{y} - \theta)^2 + vs^2 + \nu_0 s_0^2 / (n + \nu_0)$ .

When actual data  $y_d$  becomes available then conditional on the acceptance of this model inferences about  $\sigma^2$  must be made from the posterior distribution

$$p(\sigma^2 | y_d, n)(\sigma^2)^{-\left[ \frac{n + \nu_0}{2} + 1 \right]} \left( \frac{\nu_0}{\bar{\sigma}_d^2} \right)^{\frac{n + \nu_0}{2}} \exp \left[ -\frac{1}{2} \frac{(\nu_0 + n)\bar{\sigma}_d^2}{\sigma^2} \right] \quad (A.3)$$

However rational acceptance of the relevance of this model for the situation in which

$y_d$  is generated requires that relevant aspects of  $y_d$  are not surprising when assessed against a reference distribution derived from the predictive distribution.

$$p(y | n)(\sigma^2)^{-\frac{\nu_0}{2} - \left[ \frac{n + \nu_0}{2} + 1 \right]} \exp \left[ -\frac{1}{2} \frac{(\nu_0 + n)\bar{\sigma}_d^2}{\sigma^2} \right] \quad (A.4)$$



Pertinent features of the sample are  $\bar{y}_d = \bar{y}/n$ ,  $s_d^2 = 1/n \sum (y_i - \bar{y})^2$  and functions of  $(n-2)$  residual quantities  $u_1, u_2, \dots, u_{n-2}$  defined as before. These must be considered against their relevant reference distributions derived from  $p(y|M)$ . The Jacobian of the transformation from  $y$  to  $\bar{y}, s^2, \underline{u}$  is proportional to

$$(s_d^2)^{\frac{v}{2}-1} \prod_{j=1}^{n-2} \left(1 + \frac{u_j^2}{j}\right)^{-\frac{1}{2}(j+1)}$$

Thus

$$p(\bar{y}, s^2, \underline{u}|M) = p(\bar{y}|s^2, M) p(s^2|M) p(\underline{u}|M) \quad (A.6)$$

and

$$p(\bar{y}|\underline{s}^2, M) \propto \frac{1}{s_p^2} \left\{ 1 + \frac{n(\bar{y} - \theta)^2}{v s_p^2} \right\}^{-\frac{(v+1)}{2}} \quad \text{where } s_p^2 = (v s^2 + v_0 s_0^2)/(v + v_0) \quad (A.7)$$

and  $v_p = (v + v_0)$

$$p(s^2|M) \propto \frac{1}{s_0^2} \frac{1}{F^{\frac{v}{2}-1}} \frac{1}{F^{\frac{v}{2}-1}} \quad \text{where } F = \frac{s^2}{s_0^2}$$

$$p(\underline{u}|M) \propto \prod_{j=1}^{n-2} \left(1 + \frac{u_j^2}{j}\right)^{-\frac{1}{2}(j+1)} \quad (A.8)$$

Unusual features of  $s^2, \bar{y}$  and  $u_1, \dots, u_{n-2}$  would thus be assessed by computing

- (i)  $\Pr\{p(s^2|M) < p(s_d^2|M)\}$
- (ii)  $\Pr\{p(\bar{y}|\underline{s}^2, M) < p(\bar{y}_d|s_d^2, M)\}$
- (iii)  $\Pr\{p(u_j|M) < p(u_{jd}|M)\}$

These are two tailed probabilities associated with reference of

- (i)  $s^2/s_0^2$  to an  $F$  distribution with  $v$  and  $v_0$  degrees of freedom

- (ii)  $\sqrt{n}(\bar{y} - \theta)/s_p$  to a  $t^*$  table.

- (iii)  $g_j$  to the reference distribution obtained by appropriate graduation of  $p(u)$ .

Inferences about the variance

- (a) Suppose  $v_0 \rightarrow 0$ .

This limit corresponds to usual noninformative Jeffereys' prior. Again the values of  $v_0$  that could represent real situations could approach zero but not reach it.

Since in practice there could always be values of  $s^2$  which would be surprising even though  $p(\bar{y}|\underline{s}^2, M)$  was disperse, this would correspond to the situation where a very small value of  $p(F|v, v_0)$  was found even though  $v_0$  was very small.

- (b) Suppose  $v_0$  is very large

Then  $s_0^2$  and  $s_p^2 = (v s^2 + v_0 s_0^2)/(v + v_0)$  are very precisely known and if we believe the model the posterior distribution  $p(\bar{y}|\underline{s}^2, M)$ , is sharply concentrated about  $\bar{y}_0$  and our belief a posteriori is the same as that a priori. However for  $p(\bar{y}|M)$  we obtain

$$p\left(\frac{\bar{y} - \theta}{s_p/\sqrt{n}} | M\right) = p\left(z = \frac{\bar{y} - \theta}{s_p/\sqrt{n}}\right) \quad (A.10)$$

where  $z$  is a unit normal deviate and

$$p\left(\frac{s^2}{s_0^2} | M\right) = p\left\{\frac{\chi_v^2}{v} = \frac{s^2}{s_0^2}\right\} \quad (A.11)$$

So that it is only after applying the checks  $c(\bar{y})$  and  $c(s^2)$  as well as  $c_j(u)$  that we could logically use Bayes theorem.

# REFERENCES

- [1] Anscombe, F. J. and Tukey, J. W. (1963), The examination and analysis of residuals. Technometrics, 5, p.141.
- [2] Box, G.E.P. (1976), Science and Statistics, JASA, 71, p.791-799.
- [3] Box, G.E.P. and Cox, D. R. (1964), An analysis of transformations. JRSS, Series B, 26, p.211.
- [4] Box, G.E.P. and Jenkins (1976), Time Series Analysis: Forecast and Control. Holden-Day.
- [5] Box, G.E.P. and Tiao, G. C. (1962), A further look at robustness via Bayes' theorem. Biometrika, 49, p.419.
- [6] Box, G.E.P. and Tiao, G. C. (1968), A Bayesian approach to some outlier problems, Biometrika, 55, p.119.
- [7] Box, G.E.P. and Tiao, G. C. (1973), Bayesian Inference in Statistical Analysis. Addison-Wesley.
- [8] Coen, P. J., Gomme, E. D. and Kendall, M. G. (1969), Lagged Relationships in Economic Forecasting. JRSS, Series A, 132, p.133.
- [9] David, H. A., Hartley, H. O. and Pearson, E. S. (1954), The distribution of the ratio, in a single normal sample, of range to standard deviation. Biometrika, 41, p.482.
- [10] Dixon, W. J. (1950), Analysis of extreme values. Ann. Math. Statist., 21, p.27.
- [11] Durbin, J. and Watson, G. S. (1950), Testing for serial correlation in least square regression I. Biometrika, 37, p.409.
- [12] Durbin, J. and Watson, G. S. (1951), Testing for serial correlation in least square regression II. Biometrika, 38, p.159.
- [13] Ferguson, T. S. (1961), On the rejection of outliers. Proceedings of the Fourth Berkeley Symposium, 1, p.253.
- [14] Feller, W. (1968), An Introduction to Probability Theory and its Applications. Vol. 1, Wiley.
- [15] Grubbs, F. E. (1950), Sample criteria for testing outlying observations. Ann. Math. Statist., 21, p.27.

- [16] Lindley, D. V. and Smith, A.F.N. (1972), Bayes' Estimates for the Linear Model, (w/ discussion). JRSS, B, 34, p.1-41.
- [17] Palleen, L. C. (1977), Studies in the analysis of serially dependent data. Ph.D. thesis, University of Wisconsin.
- [18] Stein, C. (1956), Inadmissibility of the usual estimator for the mean of a multivariate normal distribution. Proceedings of the Third Berkeley Symposium, 1, p.197.
- [19] Tiao, G. C. and Ali, M. M. (1971), Analysis of correlated random effects: linear model with two random components, Biometrika, 58, p.37.
- [20] Tukey, J. W. (1949), One degree of freedom for non-additivity, Biometrics, 5, p.232.
- [21] Zellner, A. and Tiao, G. C. (1964), Bayesian analysis of the regression model with autocorrelated errors. JASA, 59, p.763.
- [22] Box, G.E.P. and Tiao, G. C. (1968), Bayesian analysis of means for the random effect model, J. Amer. Statist. Assoc., 63, p.179.
- [23] Box, G.E.P. and Newbold, Paul (1971), Some comments on a paper of Coen, Gomme, and Kendall, JRSS, A, 134, p.229.
- [24] Box, G.E.P. (1979), Robustness in the strategy of scientific model building, Proceedings of A.R.O. workshop at Durham, N.C., April 1978, on Robustness in Statistics, Academic Press (to appear).

1. REPORT NUMBER 1969		2. GOVT ACCESSION NO.		3. RESEARCHER'S CATALOG NUMBER	
4. TITLE (and Subtitle) SAMPLING AND BAYES' INFERENCE IN THE ADVANCEMENT OF LEARNING				5. TYPE OF REPORT & PERIOD COVERED Summary Report - no specific reporting period	
7. AUTHOR(S) George E. P. Box				8. CONTRACT OR GRANT NUMBER(S) DAG29-75-C-0024	
9. PERFORMING ORGANIZATION NAME AND ADDRESS Mathematics Research Center, University of 610 Walnut Street Madison, Wisconsin 53706				10. PROGRAM ELEMENT PROJECT, TASK AND MONITORING NUMBERS Work Unit Number 4 - Probability, Statistics, and Combinatorics	
11. CONTROLLING OFFICE NAME AND ADDRESS U. S. Army Research Office P.O. Box 12211 Research Triangle Park, North Carolina 27709				12. REPORT DATE June 1979	
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)				13. NUMBER OF PAGES 28	
				15. SECURITY CLASS. (of this report) UNCLASSIFIED	
16. DISTRIBUTION STATEMENT (for this report) Approved for public release; distribution unlimited.				15a. DECLASSIFICATION/DOWNGRADING SCHEDULE	
17. DISTRIBUTION STATEMENT (for the abstract entered in Block 20, if different from Report)					
18. SUPPLEMENTARY NOTES					
19. KEY WORDS (Continue on reverse side if necessary and identify by Block number) Bayesian inference, Sampling theory inference, Estimation, Models, Predictive distribution					
20. ABSTRACT (Continue on reverse side if necessary and identify by Block number) Scientific method is a process of guided learning in which accelerated acquisition of knowledge relevant to some question under investigation is achieved by a hierarchy of iterations in which induction and deduction are used in alternation. This process employs a developing model (or series of models implicit or explicit) against which data can be viewed. Ideally at any given stage (continued)					

DD FORM 1 JAN 73 1473 EDITION OF 1 NOV 65 IS OBSOLETE UNCLASSIFIED  
SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

of an investigation, such a model approximates relevant aspects of the underlying system and motivates the acquisition of further data as well as its analysis. By the use of a prior distribution it is possible to represent some aspects of such a model as completely known and others as more or less unknown.

Now parsimony requires that, at any given stage, the model is no more complex than is necessary to achieve a desirable degree of approximation and since each investigation is unique we cannot be sure in advance that any model we postulate will meet this goal. Therefore, at the various points in our investigation where data analysis is required, two types of inference are involved: model criticism and parameter estimation. To effect the latter, conditional on the plausibility of the model, and given the data, we can, using Bayes' Theorem, deduce posterior distributions for unknown parameters and so make inferences about them. But, before we can rely on such conditional deduction, we ought logically to check whether the model postulated accords with the data at all and, if not, consider how it should be modified. In practice, this question is usually investigated by inspecting residuals, by other informal techniques, and sometimes by making formal tests of goodness of fit. In any case model criticism, the inferential procedure whereby the need for model modification is induced, is ultimately dependent on sampling theory argument. These principles are formalized by an appropriate analysis of Bayes' formula, and implications for robust estimation are considered.