

# BAYESIAN STATISTICS

**José M. Bernardo**

*Departamento de Estadística, Facultad de Matemáticas, 46100–Burjassot, Valencia, Spain.*

**Keywords:** Amount of Information, Axiomatics, Bayes Estimator, Bayes Factor, Bayes Theorem, Credible Region, Degree of Belief, Decision Theory, Exchangeability, Foundations of Statistics, Hierarchical Models, Hypothesis Testing, Interval Estimation, Intrinsic Discrepancy, Likelihood Principle, Lindley’s Paradox, Logarithmic Divergence, Maximum Entropy, Model Choice, Model Criticism, Noninformative Priors, Nuisance Parameters, Objectivity, Point Estimation, Posterior Distribution, Predictive Distribution, Prior Choice, Prior Distribution, Probability Assessment, Probability Model, Probability Theory, Reference Distributions, Representation Theorem, Scientific Reporting, Sensitivity Analysis, Steins’s Paradox, Sufficiency.

## Contents

1. Introduction
2. Foundations
  - 2.1. Probability as a Conditional Measure of Uncertainty
  - 2.2. Statistical Inference and Decision Theory
  - 2.3. Exchangeability and Representation Theorem
3. The Bayesian Paradigm
  - 3.1. The Learning Process
  - 3.2. Predictive Distributions
  - 3.3. Asymptotic Behaviour
4. Inference Summaries
  - 4.1. Estimation
  - 4.2. Hypothesis Testing
5. Reference Analysis
  - 5.1. Reference Distributions
  - 5.2. Frequentist Properties
6. A Stylized Case Study
  - 6.1. Objective Bayesian Analysis
  - 6.2. Sensitivity Analysis
7. Discussion and Further Issues
  - 7.1. Coherency
  - 7.2. Objectivity
  - 7.3. Applicability

---

To appear at *the Encyclopedia of Life Support Systems (EOLSS)*, UNESCO. Type set on November 20, 2001.

## Glossary

**Bayes Estimator:** A function  $\omega^* = \omega^*(D)$  of data  $D$ , to be used as a proxy for the unknown value of the parameter vector  $\omega$ . It is obtained by minimizing the posterior expectation of a *loss function*,  $L(\tilde{\omega}, \omega)$ , defined to measure the consequences of using  $\tilde{\omega}$  as a proxy for the true value of  $\omega$ .

**Bayes Factor:** Given data  $D$  generated by the *probability model*  $\{p(D | \omega), \omega \in \Omega\}$  and a *prior distribution*  $p(\omega)$ , the *Bayes Factor*  $B_{01} = B_{01}(D)$  for  $\Omega_0 \subset \Omega$  against  $\Omega_1 \subset \Omega$ , is the integrated likelihood ratio  $p(D | \Omega_0)/p(D | \Omega_1)$ , where  $p(D | \Omega_i) = \int_{\Omega_i} p(D | \omega)p(\omega)d\omega$ .

**Bayes' Theorem:** Given data  $D$  generated by the *probability model*  $\{p(D | \omega), \omega \in \Omega\}$  and a *prior distribution*  $p(\omega)$ , the *posterior distribution* of  $\omega$  is  $p(\omega | D) \propto p(D | \omega)p(\omega)$ . The proportionality constant is  $\{\int_{\Omega} p(D | \omega)p(\omega)d\omega\}^{-1}$ .

**Credible Region:** Given data  $D$ , a posterior *q-credible region* for  $\omega \in \Omega$  is a subset  $R$  of  $\Omega$  with posterior probability  $q$ , so that  $\int_R p(\omega | D)d\omega = q$ . Likewise, a posterior *q-credible region* for  $\mathbf{x} \in X$  is a subset  $R$  of  $X$  with posterior predictive probability  $q$ , so that  $\int_R p(\mathbf{x} | D)d\mathbf{x} = q$ .

**Exchangeability:** The random vectors  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  are *exchangeable* if their joint distribution is invariant under permutations. An infinite sequence  $\{\mathbf{x}_j\}$  of random vectors is *exchangeable* if all its finite subsequences are exchangeable. If  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  is a *random sample* from some probability model, and hence the  $\mathbf{x}_j$ 's are independent given the model parameter, then the random vectors  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  are necessarily exchangeable.

**Expected Information:** The *expected information* from data  $D \in \mathcal{D}$  generated by a *probability model*  $\{p(D | \omega), \omega \in \Omega\}$ , about a function  $\theta = \theta(\omega)$  of the parameter vector  $\omega$ , is a functional of the *prior distribution*  $p(\omega)$ , denoted as  $I^\theta\{\mathcal{D}, p(\omega)\}$ . It is defined as the expected *logarithmic divergence*  $E_D[\int_{\Theta} p(\theta | D) \log\{p(\theta | D)/p(\theta)\}d\theta]$  of the marginal prior of  $\theta$ ,  $p(\theta)$ , from the marginal posterior of  $\theta$ ,  $p(\theta | D)$ .

**Likelihood Function:** The probability (or probability density) of the observed data  $D$  as a function of the unknown parameter vector  $\omega$ ,  $l(\omega, D) = p(D | \omega)$ .

**Logarithmic Divergence:** The *logarithmic divergence* of a probability density  $\hat{p}(\mathbf{x})$  for the random vector  $\mathbf{x} \in X$  from its true probability density  $p(\mathbf{x})$ , is the non-negative number  $\delta\{\hat{p}(\mathbf{x}) | p(\mathbf{x})\} = \int_X p(\mathbf{x}) \log\{p(\mathbf{x})/\hat{p}(\mathbf{x})\} d\mathbf{x}$ .

**Intrinsic Discrepancy Function:** Given the *probability model*  $\{p(D | \omega), \omega \in \Omega\}$ , the *intrinsic discrepancy* between the parameter values  $\omega_1$  and  $\omega_2$  is the minimum logarithmic divergence between the models  $p(D | \omega_1)$  and  $p(D | \omega_2)$ , *i.e.*, the symmetric, non-negative, function  $d(\omega_1, \omega_2) = \min(\delta\{p(D | \omega_1) | p(D | \omega_2)\}, \delta\{p(D | \omega_2) | p(D | \omega_1)\})$ .

**Maximum Likelihood Estimator (MLE):** Given data  $D$ , the maximum likelihood estimator of  $\omega \in \Omega$  is that value  $\hat{\omega} \in \Omega$  which maximizes the *likelihood function*  $l(\omega, D)$ .

**Outcome space:** See *Probability Model*.

**Parameter, Parameter Space:** See *Representation Theorem*.

**Posterior Distribution:** A probability distribution on the unknown parameter vector  $\omega \in \Omega$  in the *probability model*, typically described by its density function  $p(\omega | D)$ , which conditional on the model, encapsulates the available information about the unknown value of  $\omega$ , given the observed data  $D$  and the knowledge about  $\omega$  which the *prior distribution*  $p(\omega)$  might contain. It is obtained by *Bayes' theorem*.

**Predictive Distribution:** If  $\{\mathbf{x}_j\}$  is a sequence of exchangeable random vectors, the predictive density of a future element  $\mathbf{x}$  of the sequence, given  $n$  observed values  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , is  $p(\mathbf{x} | \mathbf{x}_1, \dots, \mathbf{x}_n) = p(\mathbf{x}, \mathbf{x}_1, \dots, \mathbf{x}_n) / p(\mathbf{x}_1, \dots, \mathbf{x}_n)$ . In particular, given a *probability model* of the form  $\{p(D | \omega), \omega \in \Omega\}$ , with  $p(D | \omega) = \prod_j p(\mathbf{x}_j | \omega)$ , and a *prior distribution*  $p(\omega)$ , the *posterior predictive density* is  $p(\mathbf{x} | D) = \int_{\Omega} p(\mathbf{x} | \omega) p(\omega | D) d\omega$ . If no data are available, the *prior predictive density* is  $p(\mathbf{x}) = \int_{\Omega} p(\mathbf{x} | \omega) p(\omega) d\omega$ .

**Prior Distribution:** A probability distribution on the unknown parameter vector  $\omega \in \Omega$  in the *probability model*, typically described by its density function  $p(\omega)$ , with  $p(\omega) \geq 0$ ,  $\int_{\Omega} p(\omega) d\omega = 1$ , which encapsulates the available information about the unknown value of  $\omega$ . If no prior information is to be utilized, the *prior distribution*  $p(\omega)$  may be replaced by a *reference prior function*  $\pi(\omega)$ .

**Probability Model:** A family of probability distributions of  $D \in \mathcal{D}$ , typically described by their density functions  $\{p(D | \omega), \omega \in \Omega\}$ , with  $p(D | \omega) \geq 0$  and  $\int_{\mathcal{D}} p(D | \omega) dD = 1$  for all  $\omega \in \Omega$ , which is assumed to contain the probability mechanism which has generated the observed data  $D$ . The set  $\mathcal{D}$  of possible data values is the *outcome space*. If data are a *random sample*  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  of observations  $\mathbf{x}_i \in X$  generated from  $p(\mathbf{x} | \omega)$ , then  $p(D | \omega) = \prod_j p(\mathbf{x}_j | \omega)$  and  $\mathcal{D} = X^n$ .

**Random Sample:** See *Probability Model*.

**Reference Prior Function:** Given the probability model  $\{p(D | \omega), \omega \in \Omega\}$  and a function  $\phi = \phi(\omega)$  of the parameter vector  $\omega$ , a  $\phi$ -reference prior is a positive function  $\pi_{\phi}(\omega)$  to be used as the prior distribution in Bayes theorem to obtain a *reference posterior distribution* for  $\phi$ . The function  $\pi_{\phi}(\omega)$  does not necessarily have a finite integral over  $\Omega$  and, hence, it is not necessarily a probability distribution.

**Reference Posterior Distribution:** Given data  $D$  and a *probability model*  $\{p(D | \omega), \omega \in \Omega\}$ , which is assumed to have generated  $D$ , the reference posterior of a function  $\phi = \phi(\omega)$  of the parameter vector  $\omega$  is a probability distribution, typically described by its density function  $\pi(\phi | D)$ , which encapsulates inferential conclusions on the value of  $\phi$ , *solely* based on the assumed model and the observed data  $D$ . It is obtained from Bayes' theorem with a  $\phi$ -reference prior function as a formal prior.

**Representation Theorem:** If  $\{\mathbf{x}_j\}$  is an infinite sequence of exchangeable random vectors, the joint density of any finite subsequence  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  has an integral representation of the form  $p(\mathbf{x}_1, \dots, \mathbf{x}_n) = \int_{\Omega} \prod_{i=1}^n p(\mathbf{x}_i | \omega) p(\omega) d\omega$ . Thus, observations  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  may be treated as a *random sample* from some distribution  $p(\mathbf{x} | \omega)$ , labeled by a *parameter*  $\omega \in \Omega$ , which is defined as the limit (as  $n \rightarrow \infty$ ) of some function of  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , and *there exists* a probability distribution  $p(\omega)$  over the *parameter space*  $\Omega$ .

**Sufficiency:** Given the *probability model*  $\{p(D | \omega), \omega \in \Omega\}$ , a function of the data  $\mathbf{t} = \mathbf{t}(D)$ , is a sufficient statistic if (and only if) there exist functions  $f$  and  $g$  such that, for any data set  $D \in \mathcal{D}$ , the *likelihood function* factorizes in the form  $l(\omega, D) = p(D | \omega) = f(\omega, \mathbf{t})g(D)$ . A necessary and sufficient condition for  $\mathbf{t}$  to be sufficient is that, for any prior  $p(\omega)$ , the posterior distribution  $p(\omega | D) = p(\omega | \mathbf{t})$  only depends on the data through the function  $\mathbf{t} = \mathbf{t}(D)$ .

## Summary

Statistics is the study of uncertainty. Bayesian statistical methods provide a *complete* paradigm for both statistical inference and decision making under uncertainty. Bayesian methods are firmly based on strict mathematical foundations, providing a *coherent* methodology which makes it possible to incorporate relevant initial information, and which solves many of the difficulties faced by conventional statistical methods. The Bayesian paradigm is based on an interpretation of probability as a *conditional measure of uncertainty* which closely matches the use of the word ‘probability’ in ordinary language. Statistical inference about a quantity of interest is described as the modification of the uncertainty about its value in the light of evidence, and Bayes’ theorem specifies how this modification should be made. Bayesian methods may be applied to complex, richly structured problems, which have been fairly inaccessible to traditional statistical methods. The special situation, often met in scientific reporting and public decision making, where the only acceptable information is that which may be deduced from available documented data, is addressed as an important particular case.

## 1. Introduction

Scientific experimental or observational results generally consist of (possibly many) sets of data of the general form  $D = \{x_1, \dots, x_n\}$ , where the  $x_i$ ’s are somewhat “homogeneous” (possibly multidimensional) observations  $x_i$ . Statistical methods are then typically used to derive conclusions on both the nature of the process which has produced those observations, and on the expected behaviour of future instances of the same process. A central element of *any* statistical analysis is the specification of a *probability model* which is assumed to describe the mechanism which has generated the observed data  $D$  as a function of a (possibly multidimensional) parameter  $\omega \in \Omega$ , sometimes named the *state of nature*, about whose value only limited information (if any) is available. All derived statistical conclusions are obviously conditional on the assumed probability model.

Unlike most other branches of mathematics, conventional methods of statistical inference suffer from the lack of an axiomatic basis; as a consequence, their proposed desiderata are often mutually incompatible, and the analysis of the same data may well lead to incompatible results when different, apparently intuitive procedures are tried. In marked contrast, the Bayesian approach to statistical inference is firmly based on axiomatic foundations which provide a unifying logical structure, and guarantee the mutual consistency of the methods proposed. Bayesian methods constitute a *complete* paradigm to statistical inference, a scientific revolution in Kuhn’s sense.

Bayesian statistics only require the *mathematics* of probability theory and the *interpretation* of probability which most closely corresponds to the standard use of this word in everyday language: it is no accident that some of the more important seminal books on Bayesian statistics, such as the works of de Laplace, de Finetti or Jeffreys, are actually entitled “Probability Theory”. The practical consequences of adopting the Bayesian paradigm are far reaching. Indeed, Bayesian methods (i) reduce statistical inference to problems in probability theory, thereby minimizing the need for completely new concepts, and (ii) serve to discriminate among conventional statistical techniques, by either providing a logical justification to some (and making explicit the conditions under which they are valid), or proving the logical inconsistency of others.

The main consequence of these foundations is the mathematical *need* to describe by means of probability distributions all uncertainties present in the problem. In particular, unknown parameters in probability models *must* have a joint probability distribution which describes the available information about their values; this is often regarded as the more characteristic element

of a Bayesian approach. Notice that (in sharp contrast to conventional statistics) *parameters are treated as random variables* within the Bayesian paradigm. This is not a description of their variability (parameters are typically *fixed unknown* quantities) but a description of the *uncertainty* about their true values.

An important particular case arises when either no relevant prior information is readily available, or that information is subjective and an “objective” analysis is desired, one exclusively based on accepted model assumptions and well-documented data. This is addressed by *reference analysis* which uses information-theoretical concepts to derive appropriate reference posterior distributions, defined to encapsulate inferential conclusions on the quantities of interest solely based on the supposed model and the observed data.

In this article it is assumed that probability distributions may be described through their probability density functions, and no distinction is made between a random quantity and the particular values that it may take. Bold roman fonts are used for *observable* random vectors (typically data) and bold greek fonts are used for unobservable random vectors (typically parameters); lower case is used for variables and upper case for their dominion sets. Moreover, the standard mathematical convention of referring to *functions*, say  $f$  and  $g$  of  $\mathbf{x} \in X$ , respectively by  $f(\mathbf{x})$  and  $g(\mathbf{x})$ , will be used throughout. Thus,  $p(\boldsymbol{\theta} | C)$  and  $p(\mathbf{x} | C)$  respectively represent general *probability densities* of the random vectors  $\boldsymbol{\theta} \in \Theta$  and  $\mathbf{x} \in X$  under conditions  $C$ , so that  $p(\boldsymbol{\theta} | C) \geq 0$ ,  $\int_{\Theta} p(\boldsymbol{\theta} | C) d\boldsymbol{\theta} = 1$ , and  $p(\mathbf{x} | C) \geq 0$ ,  $\int_X p(\mathbf{x} | C) d\mathbf{x} = 1$ . This admittedly unprecise notation will greatly simplify the exposition. If the random vectors are discrete, these functions naturally become probability mass functions, and integrals over their values become sums.

**Table 1.** Notation for common probability density and probability mass functions

Name	Probability Density or Probability Mass Function	Parameter(s)
Beta	$\text{Be}(x   \alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1}(1-x)^{\beta-1}$ , $x \in (0, 1)$	$\alpha > 0, \beta > 0$
Binomial	$\text{Bi}(x   n, \theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$ , $x \in \{0, \dots, n\}$	$n \in \{1, 2, \dots\}, \theta \in (0, 1)$
Exponential	$\text{Ex}(x   \theta) = \theta e^{-\theta x}$ , $x > 0$	$\theta > 0$
ExpGamma	$\text{Eg}(x   \alpha, \beta) = \frac{\alpha\beta^\alpha}{(x+\beta)^{\alpha+1}}$ , $x > 0$	$\alpha > 0, \beta > 0$
Gamma	$\text{Ga}(x   \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}$ , $x > 0$	$\alpha > 0, \beta > 0$
NegBinomial	$\text{Nb}(x   r, \theta) = \theta^r \binom{r+x-1}{r-1} (1-\theta)^x$ , $x \in \{0, 1, \dots\}$	$r \in \{1, 2, \dots\}, \theta \in (0, 1)$
Normal	$\text{N}_k(\mathbf{x}   \boldsymbol{\mu}, \Sigma) = \frac{ \Sigma ^{-1/2}}{(2\pi)^{k/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^t \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$ , $\mathbf{x} \in \mathfrak{R}^k$	$\boldsymbol{\mu} \in \mathfrak{R}^k, \Sigma \text{ def. pos.}$
Poisson	$\text{Pn}(x   \lambda) = e^{-\lambda} \frac{\lambda^x}{x!}$ , $x \in \{0, 1, \dots\}$	$\lambda > 0$
Student	$\text{St}(x   \mu, \sigma, \alpha) = \frac{\Gamma(\frac{\alpha+1}{2})}{\Gamma(\frac{\alpha}{2})} \frac{1}{\sigma\sqrt{\alpha\pi}} \left[1 + \frac{1}{\alpha} \left(\frac{x-\mu}{\sigma}\right)^2\right]^{-(\alpha+1)/2}$ , $x \in \mathfrak{R}$	$\mu \in \mathfrak{R}, \sigma > 0, \alpha > 0$

Specific density functions are denoted by appropriate names. Thus, if  $x$  is a random quantity with a normal distribution of mean  $\mu$  and standard deviation  $\sigma$ , its probability density function will be denoted  $\text{N}(x | \mu, \sigma)$ . Table 1 contains definitions of other distributions used in this article.

Bayesian methods make frequent use of the logarithmic divergence, a very general measure of the goodness of the approximation of a probability density  $p(\mathbf{x})$  by another density  $\hat{p}(\mathbf{x})$ . The *logarithmic divergence* of a probability density  $\hat{p}(\mathbf{x})$  of the random vector  $\mathbf{x} \in X$  from its true probability density  $p(\mathbf{x})$ , is defined as  $\delta\{\hat{p}(\mathbf{x}) | p(\mathbf{x})\} = \int_X p(\mathbf{x}) \log\{p(\mathbf{x})/\hat{p}(\mathbf{x})\} d\mathbf{x}$ . It may be shown that (i) the logarithmic divergence is non-negative (and it is zero if, and only if,  $\hat{p}(\mathbf{x}) = p(\mathbf{x})$  almost everywhere), and (ii) that  $\delta\{\hat{p}(\mathbf{x}) | p(\mathbf{x})\}$  is invariant under one-to-one transformations of  $\mathbf{x}$ .

This article contains a brief summary of the mathematical foundations of Bayesian statistical methods (Section 2), an overview of the paradigm (Section 3), a description of useful inference summaries, including estimation and hypothesis testing (Section 4), an explicit discussion of objective Bayesian methods (Section 5), the detailed analysis of a simplified case study (Section 6), and a final discussion which includes pointers to further issues not addressed in the main text (Section 7).

## 2. Foundations

A central element of the Bayesian paradigm is the use of probability distributions to describe all relevant unknown quantities, interpreting the probability of an event as a conditional measure of uncertainty, on a  $[0, 1]$  scale, about the occurrence of the event in some specific conditions. The limiting extreme values 0 and 1, which are typically inaccessible in applications, respectively describe impossibility and certainty of the occurrence of the event. This interpretation of probability includes and extends all other probability interpretations. There are two independent arguments which prove the mathematical inevitability of the use of probability distributions to describe uncertainties; these are summarized later in this section.

### 2.1. Probability as a Measure of Conditional Uncertainty

Bayesian statistics uses the word *probability* in precisely the same sense in which this word is used in everyday language, as a *conditional measure of uncertainty* associated with the occurrence of a particular event, given the available information and the accepted assumptions. Thus,  $\Pr(E | C)$  is a measure of (presumably rational) belief in the occurrence of the *event*  $E$  under *conditions*  $C$ . It is important to stress that probability is *always* a function of two arguments, the event  $E$  whose uncertainty is being measured, and the conditions  $C$  under which the measurement takes place; “absolute” probabilities do not exist. In typical applications, one is interested in the probability of some event  $E$  given the available *data*  $D$ , the set of *assumptions*  $A$  which one is prepared to make about the mechanism which has generated the data, and the relevant contextual *knowledge*  $K$  which might be available. Thus,  $\Pr(E | D, A, K)$  is to be interpreted as a measure of (presumably rational) belief in the occurrence of the *event*  $E$ , given data  $D$ , assumptions  $A$  and any other available knowledge  $K$ , as a measure of how “likely” is the occurrence of  $E$  in these conditions. Sometimes, but certainly not always, the probability of an event under given conditions may be associated with the relative frequency of “similar” events in “similar” conditions. The following examples are intended to illustrate the use of probability as a conditional measure of uncertainty.

*Probabilistic diagnosis.* A human population is known to contain 0.2% of people infected by a particular virus. A person, *randomly selected* from that population, is subject to a test which is known from laboratory data to yield positive results in 98% of infected people and in 1% of non-infected, so that, if  $V$  denotes the event that a person carries the virus and  $+$  denotes a positive result,  $\Pr(+ | V) = 0.98$  and  $\Pr(+ | \bar{V}) = 0.01$ . Suppose that the result of the test turns out to be positive. Clearly, one is then interested in  $\Pr(V | +, A, K)$ , the *probability* that the person

carries the virus, given the positive result, the assumptions  $A$  about the probability mechanism generating the test results, and the available knowledge  $K$  of the prevalence of the infection in the population under study (described here by  $\Pr(V | K) = 0.002$ ). An elementary exercise in probability algebra, which involves Bayes' theorem in its simplest form (see Section 3), yields  $\Pr(V | +, A, K) = 0.164$ . Notice that the four probabilities involved in the problem have *precisely the same interpretation*: they are all conditional measures of uncertainty. Besides,  $\Pr(V | +, A, K)$  is *both* a measure of the uncertainty associated with the event that the particular person who tested positive is actually infected, *and* an *estimate* of the proportion of people in that population (about 16.4%) that would eventually prove to be infected among those which yielded a positive test. ◁

*Estimation of a proportion.* A survey is conducted to estimate the proportion  $\theta$  of individuals in a population who share a given property. A random sample of  $n$  elements is analyzed,  $r$  of which are found to possess that property. One is then typically interested in using the results from the sample to establish regions of  $[0, 1]$  where the unknown value of  $\theta$  may plausibly be expected to lie; this information is provided by *probabilities* of the form  $\Pr(a < \theta < b | r, n, A, K)$ , a conditional measure of the uncertainty about the event that  $\theta$  belongs to  $(a, b)$  given the information provided by the data  $(r, n)$ , the assumptions  $A$  made on the behaviour of the mechanism which has generated the data (a random sample of  $n$  Bernoulli trials), and any relevant knowledge  $K$  on the values of  $\theta$  which might be available. For example, after a political survey in which 720 citizens out of a random sample of 1500 have declared to be in favour of a particular political measure, one may conclude that  $\Pr(\theta < 0.5 | 720, 1500, A, K) = 0.933$ , indicating a probability of about 93% that a referendum of that issue would be lost. Similarly, after a screening test for an infection where 100 people have been tested, none of which has turned out to be infected, one may conclude that  $\Pr(\theta < 0.01 | 0, 100, A, K) = 0.844$ , or a probability of about 84% that the proportion of infected people is smaller than 1%. ◁

*Measurement of a physical constant.* A team of scientists, intending to establish the unknown value of a physical constant  $\mu$ , obtain data  $D = \{x_1, \dots, x_n\}$  which are considered to be measurements of  $\mu$  subject to error. The probabilities of interest are then typically of the form  $\Pr(a < \mu < b | x_1, \dots, x_n, A, K)$ , the *probability* that the unknown value of  $\mu$  (fixed in nature, but unknown to the scientists) lies within an interval  $(a, b)$  given the information provided by the data  $D$ , the assumptions  $A$  made on the behaviour of the measurement mechanism, and whatever knowledge  $K$  might be available on the value of the constant  $\mu$ . Again, those probabilities are conditional measures of uncertainty which describe the (necessarily probabilistic) conclusions of the scientists on the true value of  $\mu$ , given available information and accepted assumptions. For example, after a classroom experiment to measure the gravitational field with a pendulum, a student may report (in  $\text{m/sec}^2$ ) something like  $\Pr(9.788 < g < 9.829 | D, A, K) = 0.95$ , meaning that, under accepted knowledge  $K$  and assumptions  $A$ , the *observed* data  $D$  indicate that the true value of  $g$  lies within 9.788 and 9.829 with probability 0.95, a conditional uncertainty measure on a  $[0, 1]$  scale. This is naturally compatible with the fact that the value of the gravitational field at the laboratory may well be known with high precision from available literature or from precise previous experiments, but the student may have been instructed *not* to use that information as part of the accepted knowledge  $K$ . Under some conditions, it is also true that if the same *procedure* were actually used by many other students with similarly obtained data sets, their reported intervals would actually cover the true value of  $g$  in approximately 95% of the cases, thus providing some form of *calibration* for the student's probability statement (see Section 5.2). ◁

*Prediction.* An experiment is made to count the number  $r$  of times that an event  $E$  takes place in each of  $n$  replications of a well defined situation; it is observed that  $E$  does take place  $r_i$  times in replication  $i$ , and it is desired to forecast the number of times  $r$  that  $E$  will take place in a future, similar situation. This is a *prediction* problem on the value of an *observable* (discrete) quantity  $r$ , given the information provided by data  $D$ , accepted assumptions  $A$  on the probability mechanism which generates the  $r_i$ 's, and any relevant available knowledge  $K$ . Hence, simply the computation of the probabilities  $\{\Pr(r | r_1, \dots, r_n, A, K)\}$ , for  $r = 0, 1, \dots$ , is required. For example, the quality assurance engineer of a firm which produces automobile restraint systems may report something like  $\Pr(r = 0 | r_1 = \dots = r_{10} = 0, A, K) = 0.953$ , after observing that the entire production of airbags in each of  $n = 10$  consecutive months has yielded no complaints from their clients. This should be regarded as a measure, on a  $[0, 1]$  scale, of the conditional uncertainty, given observed data, accepted assumptions and contextual knowledge, associated with the event that no airbag complaint will come from next month's production and, if conditions remain constant, this is also an estimate of the proportion of months expected to share this desirable property.

A similar problem may naturally be posed with continuous observables. For instance, after measuring some continuous magnitude in each of  $n$  randomly chosen elements within a population, it may be desired to forecast the proportion of items in the whole population whose magnitude satisfies some precise specifications. As an example, after measuring the breaking strengths  $\{x_1, \dots, x_{10}\}$  of 10 randomly chosen safety belt webbings to verify whether or not they satisfy the requirement of remaining above 26 kN, the quality assurance engineer may report something like  $\Pr(x > 26 | x_1, \dots, x_{10}, A, K) = 0.9987$ . This should be regarded as a measure, on a  $[0, 1]$  scale, of the conditional uncertainty (given observed data, accepted assumptions and contextual knowledge) associated with the event that a randomly chosen safety belt webbing will support no less than 26 kN. If production conditions remain constant, it will also be an estimate of the proportion of safety belts which will conform to this particular specification.

Often, additional information of future observations is provided by related covariates. For instance, after observing the outputs  $\{y_1, \dots, y_n\}$  which correspond to a sequence  $\{x_1, \dots, x_n\}$  of different production conditions, it may be desired to forecast the output  $y$  which would correspond to a particular set  $x$  of production conditions. For instance, the viscosity of commercial condensed milk is required to be within specified values  $a$  and  $b$ ; after measuring the viscosities  $\{y_1, \dots, y_n\}$  which correspond to samples of concentrated milk produced under different physical conditions  $\{x_1, \dots, x_n\}$ , production engineers will require probabilities of the form  $\Pr(a < y < b | x, (y_1, x_1), \dots, (y_n, x_n), A, K)$ . This is a conditional measure of the uncertainty (always given observed data, accepted assumptions and contextual knowledge) associated with the event that condensed milk produced under conditions  $x$  will actually satisfy the required viscosity specifications. ◁

## 2.2. Statistical Inference and Decision Theory

Decision theory not only provides a precise methodology to deal with decision problems under uncertainty, but its solid axiomatic basis also provides a powerful case for the logical force of the Bayesian approach. We now summarize the basic argument.

A decision problem exists whenever there are two or more possible courses of action; let  $\mathcal{A}$  be the class of possible actions. Moreover, for each  $a \in \mathcal{A}$ , let  $\Theta_a$  be the set of *relevant events* which may affect the result of choosing  $a$ , and let  $c(a, \theta) \in \mathcal{C}_a$ ,  $\theta \in \Theta_a$ , be the *consequence* of having chosen action  $a$  when event  $\theta$  takes place. The class of pairs  $\{(\Theta_a, \mathcal{C}_a), a \in \mathcal{A}\}$  describes the *structure* of the decision problem. Without loss of generality, it may be assumed that the



possible actions are mutually exclusive, for otherwise one would work with the appropriate Cartesian product.

Different sets of principles have been proposed to capture a minimum collection of logical rules that could sensibly be required for “rational” decision-making. These all consist of axioms with a strong intuitive appeal; examples include the *transitivity* of preferences (if  $a_1 > a_2$  given  $C$ , and  $a_2 > a_3$  given  $C$ , then  $a_1 > a_3$  given  $C$ ), and the *sure-thing principle* (if  $a_1 > a_2$  given  $C$  and  $E$ , and  $a_1 > a_2$  given  $C$  and  $\bar{E}$ , then  $a_1 > a_2$  given  $C$ ). Notice that these rules are not intended as a description of actual human decision-making, but as a *normative* set of principles to be followed by someone who aspires to coherent decision-making.

There are naturally different options for the set of acceptable principles, but all of them lead basically to the same conclusions, namely:

(i) Preferences among consequences should be measured with a real-valued bounded *utility* function  $U(c) = U(a, \theta)$  which specifies, on some numerical scale, their desirability.

(ii) The uncertainty of relevant events should be measured with a set of *probability* distributions  $\{p(\theta | C, a), \theta \in \Theta_a, a \in \mathcal{A}\}$  describing their plausibility given the conditions  $C$  under which the decision must be taken.

(iii) The desirability of the available actions is measured by their corresponding *expected utility*

$$\bar{U}(a | C) = \int_{\Theta_a} U(a, \theta) p(\theta | C, a) d\theta, \quad a \in \mathcal{A}. \quad (1)$$

It is often convenient to work in terms of the non-negative *loss* function defined by

$$L(a, \theta) = \sup_{a \in \mathcal{A}} \{U(a, \theta)\} - U(a, \theta), \quad (2)$$

which directly measures, as a function of  $\theta$ , the “penalty” for choosing a wrong action. The relative undesirability of available actions  $a \in \mathcal{A}$  is then measured by their *expected loss*

$$\bar{L}(a | C) = \int_{\Theta_a} L(a, \theta) p(\theta | C, a) d\theta, \quad a \in \mathcal{A}. \quad (3)$$

Notice that, in particular, the argument described above establishes the need to quantify the uncertainty about all relevant unknown quantities (the actual values of the  $\theta$ 's), and specifies that this quantification *must* have the mathematical structure of probability distributions. These probabilities are conditional on the circumstances  $C$  under which the decision is to be taken, which typically, but not necessarily, include the results  $D$  of some relevant experimental or observational data.

It has been argued that the development described above (which is not questioned when decisions have to be made) does not apply to problems of statistical inference, where no specific decision making is envisaged. However, there are two powerful counterarguments to this. Indeed, (i) a problem of statistical inference is typically considered worth analyzing because it *may* eventually help to make sensible decisions (as Ramsey put it in the 1930's, a lump of arsenic is poisonous because it *may* kill someone, not because it has actually killed someone), and (ii) it has been shown (by Bernardo in the 1970's) that statistical inference on  $\theta$  actually *has* the mathematical structure of a decision problem, where the class of alternatives is the functional space

$$\mathcal{A} = \left\{ p(\theta | D); \quad p(\theta | D) > 0, \quad \int_{\Theta} p(\theta | D) d\theta = 1 \right\} \quad (4)$$

of the conditional probability distributions of  $\theta$  given the data, and the utility function is a measure of the amount of information about  $\theta$  which the data may be expected to provide.

### 2.3. Exchangeability and Representation Theorem

Available data often take the form of a set  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  of “homogeneous” observations, in the precise sense that only their *values* matter and not the *order* in which they appear. Formally, this is captured by the notion of *exchangeability*. The set of random vectors  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  is exchangeable if their joint distribution is invariant under permutations. An infinite sequence  $\{\mathbf{x}_j\}$  of random vectors is exchangeable if all its finite subsequences are exchangeable. Notice that, in particular, any random sample from any model is exchangeable in this sense. The concept of exchangeability, introduced by de Finetti in the 1930’s, is central to modern statistical thinking. Indeed, the general *representation theorem* implies that if a set of observations is assumed to be a subset of an exchangeable sequence, then it constitutes a *random sample* from some probability model  $\{p(\mathbf{x} | \boldsymbol{\omega}), \boldsymbol{\omega} \in \Omega\}$ ,  $\mathbf{x} \in X$ , labeled by some *parameter vector*  $\boldsymbol{\omega}$ ; furthermore this parameter  $\boldsymbol{\omega}$  is *defined* as the limit (as  $n \rightarrow \infty$ ) of some function of the observations. Available information about the value of  $\boldsymbol{\omega}$  in prevailing conditions  $C$  is *necessarily* described by *some* probability distribution  $p(\boldsymbol{\omega} | C)$ .

For example, in the case of a sequence  $\{x_1, x_2, \dots\}$  of dichotomous exchangeable random quantities  $x_j \in \{0, 1\}$ , de Finetti’s representation theorem establishes that the joint distribution of  $(x_1, \dots, x_n)$  has an *integral representation* of the form

$$p(x_1, \dots, x_n | C) = \int_0^1 \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i} p(\theta | C) d\theta, \quad \theta = \lim_{n \rightarrow \infty} \frac{r}{n}, \quad (5)$$

where  $r = \sum x_j$  is the number of positive trials. This is precisely the joint distribution of a set of (conditionally) independent Bernoulli trials with parameter  $\theta$ , over which some probability distribution  $p(\theta | C)$  is therefore proven to exist. More generally, for sequences of arbitrary random quantities  $\{\mathbf{x}_1, \mathbf{x}_2, \dots\}$ , exchangeability leads to integral representations of the form

$$p(\mathbf{x}_1, \dots, \mathbf{x}_n | C) = \int_{\Omega} \prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\omega}) p(\boldsymbol{\omega} | C) d\boldsymbol{\omega}, \quad (6)$$

where  $\{p(\mathbf{x} | \boldsymbol{\omega}), \boldsymbol{\omega} \in \Omega\}$  denotes some probability *model*,  $\boldsymbol{\omega}$  is the limit as  $n \rightarrow \infty$  of some function  $f(\mathbf{x}_1, \dots, \mathbf{x}_n)$  of the observations, and  $p(\boldsymbol{\omega} | C)$  is some probability distribution over  $\Omega$ . This formulation includes “nonparametric” (distribution free) modelling, where  $\boldsymbol{\omega}$  may index, for instance, all continuous probability distributions on  $X$ . Notice that  $p(\boldsymbol{\omega} | C)$  does *not* describe a possible variability of  $\boldsymbol{\omega}$  (since  $\boldsymbol{\omega}$  will typically be a fixed *unknown* vector), but a description on the uncertainty associated with its actual value.

Under appropriate conditioning, exchangeability is a very general assumption, a powerful extension of the traditional concept of a *random sample*. Indeed, many statistical analyses directly assume data (or subsets of the data) to be a random sample of conditionally independent observations from some probability model, so that  $p(\mathbf{x}_1, \dots, \mathbf{x}_n | \boldsymbol{\omega}) = \prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\omega})$ ; but *any* random sample is exchangeable, since  $\prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\omega})$  is obviously invariant under permutations. Notice that the observations in a random sample are only independent *conditional* on the parameter value  $\boldsymbol{\omega}$ ; as nicely put by Lindley, the mantra that the observations  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  in a random sample are independent is ridiculous when they are used to infer  $\mathbf{x}_{n+1}$ . Notice also that, under exchangeability, the general representation theorem provides an *existence theorem* for a probability distribution  $p(\boldsymbol{\omega} | C)$  on the parameter space  $\Omega$ , and that this is an argument which only depends on mathematical probability theory.

Another important consequence of exchangeability is that it provides a formal *definition* of the parameter  $\boldsymbol{\omega}$  which labels the model as the limit, as  $n \rightarrow \infty$ , of *some* function  $f(\mathbf{x}_1, \dots, \mathbf{x}_n)$  of the observations; the function  $f$  obviously depends both on the assumed model and the chosen

parametrization. For instance, in the case of a sequence of Bernoulli trials, the parameter  $\theta$  is *defined* as the limit, as  $n \rightarrow \infty$ , of the relative frequency  $r/n$ . It follows that, under exchangeability, the sentence “the true value of  $\omega$ ” has a well-defined meaning, if only asymptotically verifiable. Moreover, if two different models have parameters which are functionally related by their definition, then the corresponding posterior distributions may be meaningfully compared, for they refer to functionally related quantities. For instance, if a finite subset  $\{x_1, \dots, x_n\}$  of an exchangeable sequence of integer observations is assumed to be a random sample from a Poisson distribution  $\text{Po}(x | \lambda)$ , so that  $E[x | \lambda] = \lambda$ , then  $\lambda$  is *defined* as  $\lim_{n \rightarrow \infty} \{\bar{x}_n\}$ , where  $\bar{x}_n = \sum_j x_j/n$ ; similarly, if for some fixed non-zero integer  $r$ , the same data are assumed to be a random sample for a negative binomial  $\text{Nb}(x | r, \theta)$ , so that  $E[x | \theta, r] = r(1 - \theta)/\theta$ , then  $\theta$  is *defined* as  $\lim_{n \rightarrow \infty} \{r/(\bar{x}_n + r)\}$ . It follows that  $\theta \equiv r/(\lambda + r)$  and, hence,  $\theta$  and  $r/(\lambda + r)$  may be treated as the *same* (unknown) quantity whenever this might be needed as, for example, when comparing the relative merits of these alternative probability models.

### 3. The Bayesian Paradigm

The statistical analysis of some observed data  $D$  typically begins with some informal *descriptive* evaluation, which is used to suggest a tentative, formal *probability model*  $\{p(D | \omega), \omega \in \Omega\}$  assumed to represent, for some (unknown) value of  $\omega$ , the probabilistic mechanism which has generated the observed data  $D$ . The arguments outlined in Section 2 establish the logical need to assess a *prior* probability distribution  $p(\omega | K)$  over the parameter space  $\Omega$ , describing the available knowledge  $K$  about the value of  $\omega$  prior to the data being observed. It then follows from standard probability theory that, if the probability model is correct, all available information about the value of  $\omega$  after the data  $D$  have been observed is contained in the corresponding *posterior* distribution whose probability density,  $p(\omega | D, A, K)$ , is immediately obtained from Bayes’ theorem,

$$p(\omega | D, A, K) = \frac{p(D | \omega) p(\omega | K)}{\int_{\Omega} p(D | \omega) p(\omega | K) d\omega}, \quad (7)$$

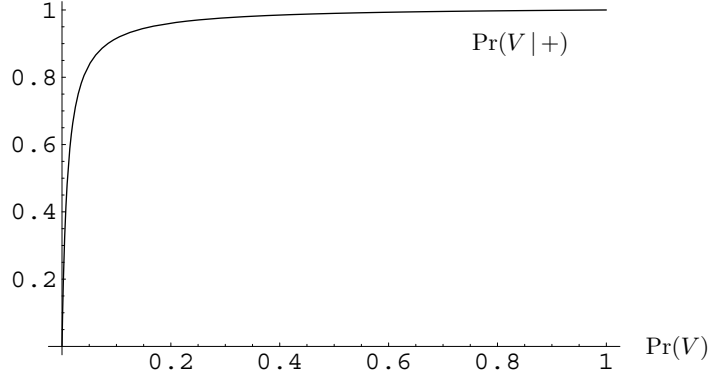
where  $A$  stands for the assumptions made on the probability model. It is this systematic use of Bayes’ theorem to incorporate the information provided by the data that justifies the adjective *Bayesian* by which the paradigm is usually known. It is obvious from Bayes’ theorem that any value of  $\omega$  with zero prior density will have zero posterior density. Thus, it is typically assumed (by appropriate restriction, if necessary, of the *parameter space*  $\Omega$ ) that prior distributions are *strictly positive* (as Savage put it, keep the mind open, or at least ajar). To simplify the presentation, the accepted assumptions  $A$  and the available knowledge  $K$  are often omitted from the notation, but the fact that *all* statements about  $\omega$  given  $D$  are *also* conditional to  $A$  and  $K$  should always be kept in mind.

*Example 1. (Bayesian inference with a finite parameter space).* Let  $p(D | \theta), \theta \in \{\theta_1, \dots, \theta_m\}$ , be the probability mechanism which is assumed to have generated the observed data  $D$ , so that  $\theta$  may only take a *finite* number of values. Using the finite form of Bayes’ theorem, and omitting the prevailing conditions from the notation, the posterior probability of  $\theta_i$  after data  $D$  have been observed is

$$\text{Pr}(\theta_i | D) = \frac{p(D | \theta_i) \text{Pr}(\theta_i)}{\sum_{j=1}^m p(D | \theta_j) \text{Pr}(\theta_j)}, \quad i = 1, \dots, m. \quad (8)$$

For any prior distribution  $p(\theta) = \{\text{Pr}(\theta_1), \dots, \text{Pr}(\theta_m)\}$  describing available knowledge about the value of  $\theta$ ,  $\text{Pr}(\theta_i | D)$  measures how likely should  $\theta_i$  be judged, given both the initial knowledge described by the prior distribution, and the information provided by the data  $D$ .

An important, frequent application of this simple technique is provided by probabilistic diagnosis. For example, consider the simple situation where a particular test designed to detect a virus is known from laboratory research to give a positive result in 98% of infected people and in 1% of non-infected. Then, the posterior probability that a person who tested positive is infected is given by  $\Pr(V | +) = (0.98 p) / \{0.98 p + 0.01 (1 - p)\}$  as a function of  $p = \Pr(V)$ , the prior probability of a person being infected (the *prevalence* of the infection in the population under study). Figure 1 shows  $\Pr(V | +)$  as a function of  $\Pr(V)$ .



**Figure 1.** Posterior probability of infection  $\Pr(V | +)$  given a positive test, as a function of the prior probability of infection  $\Pr(V)$ .

As one would expect, the posterior probability is only zero if the prior probability is zero (so that it is *known* that the population is free of infection) and it is only one if the prior probability is one (so that it is *known* that the population is universally infected). Notice that if the infection is rare, then the posterior probability of a randomly chosen person being infected will be relatively low even if the test is positive. Indeed, for say  $\Pr(V) = 0.002$ , one finds  $\Pr(V | +) = 0.164$ , so that in a population where only 0.2% of individuals are infected, only 16.4% of those testing positive within a random sample will actually prove to be infected: most positives would actually be *false positives*. ◁

In this section, we describe in some detail the learning process described by Bayes' theorem, discuss its implementation in the presence of nuisance parameters, show how it can be used to forecast the value of future observations, and analyze its large sample behaviour.

### 3.1. The Learning Process

In the Bayesian paradigm, the process of learning from the data is systematically implemented by making use of Bayes' theorem to combine the available prior information with the information provided by the data to produce the required posterior distribution. Computation of posterior densities is often facilitated by noting that Bayes' theorem may be simply expressed as

$$p(\omega | D) \propto p(D | \omega) p(\omega), \tag{9}$$

(where  $\propto$  stands for 'proportional to' and where, for simplicity, the accepted assumptions  $A$  and the available knowledge  $K$  have been omitted from the notation), since the missing proportionality constant  $[\int_{\Omega} p(D | \omega) p(\omega) d\omega]^{-1}$  may always be deduced from the fact that  $p(\omega | D)$ , a probability density, must integrate to one. Hence, to identify the form of a posterior distribution it suffices to identify a *kernel* of the corresponding probability density, that is a function  $k(\omega)$  such that  $p(\omega | D) = c(D) k(\omega)$  for some  $c(D)$  which does not involve  $\omega$ . In the examples which follow, this technique will often be used.

An *improper prior function* is defined as a positive function  $\pi(\omega)$  such that  $\int_{\Omega} \pi(\omega) d\omega$  is not finite. Equation (9), the formal expression of Bayes' theorem, remains technically valid if  $p(\omega)$  is replaced by an improper prior function  $\pi(\omega)$  provided the proportionality constant exists, thus leading to a well defined *proper* posterior density  $\pi(\omega | D) \propto p(D | \omega)\pi(\omega)$ . It will later be established (Section 5) that Bayes' theorem also remains philosophically valid if  $p(\omega)$  is replaced by an appropriately chosen reference “noninformative” (typically improper) prior function  $\pi(\omega)$ .

Considered as a function of  $\omega$ ,  $l(\omega, D) = p(D | \omega)$  is often referred to as the *likelihood function*. Thus, Bayes' theorem is simply expressed in words by the statement that *the posterior is proportional to the likelihood times the prior*. It follows from equation (9) that, provided the *same* prior  $p(\omega)$  is used, two different data sets  $D_1$  and  $D_2$ , with possibly different probability models  $p_1(D_1 | \omega)$  and  $p_2(D_2 | \omega)$  but yielding *proportional* likelihood functions, will produce identical posterior distributions for  $\omega$ . This immediate consequence of Bayes theorem has been proposed as an independent principle, the *likelihood principle*, and it is seen by many as an obvious requirement for reasonable statistical inference. In particular, for any given prior  $p(\omega)$ , the posterior distribution does not depend on the set of possible data values, or *outcome space*. Notice, however, that the likelihood principle only applies to inferences about the parameter vector  $\omega$  once the data have been obtained. Consideration of the outcome space is essential, for instance, in model criticism, in the design of experiments, in the derivation of predictive distributions, or (see Section 5) in the construction of objective Bayesian procedures.

Naturally, the terms prior and posterior are only *relative* to a particular set of data. As one would expect from the coherence induced by probability theory, if data  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  are sequentially presented, the final result will be the same whether data are globally or sequentially processed. Indeed,  $p(\omega | \mathbf{x}_1, \dots, \mathbf{x}_{i+1}) \propto p(\mathbf{x}_{i+1} | \omega) p(\omega | \mathbf{x}_1, \dots, \mathbf{x}_i)$ , for  $i = 1, \dots, n - 1$ , so that the “posterior” at a given stage becomes the “prior” at the next.

In most situations, the posterior distribution is “narrower” than the prior so that, in most cases,  $p(\omega | \mathbf{x}_1, \dots, \mathbf{x}_{i+1})$  will be more concentrated around the true value of  $\omega$  than  $p(\omega | \mathbf{x}_1, \dots, \mathbf{x}_i)$ . However, this is not always the case: occasionally, a “surprising” observation will increase, rather than decrease, the uncertainty about the value of  $\omega$ . For instance, in probabilistic diagnosis, a sharp posterior probability distribution (over the possible causes  $\{\omega_1, \dots, \omega_k\}$  of a syndrome) describing, a “clear” diagnosis of disease  $\omega_i$  (that is, a posterior with a large probability for  $\omega_i$ ) would typically update to a less concentrated posterior probability distribution over  $\{\omega_1, \dots, \omega_k\}$  if a new clinical analysis yielded data which were unlikely under  $\omega_i$ .

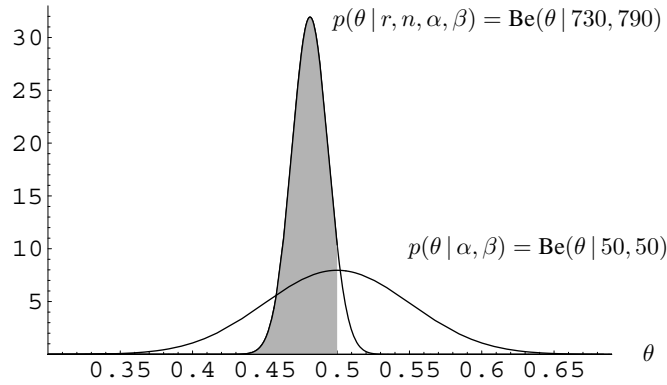
For a given probability model, one may find that a particular function of the data  $\mathbf{t} = \mathbf{t}(D)$  is a *sufficient* statistic in the sense that, given the model,  $\mathbf{t}(D)$  contains all information about  $\omega$  which is available in  $D$ . Formally,  $\mathbf{t} = \mathbf{t}(D)$  is sufficient if (and only if) there exist nonnegative functions  $f$  and  $g$  such that the likelihood function may be factorized in the form  $p(D | \omega) = f(\omega, \mathbf{t})g(D)$ . A sufficient statistic always exists, for  $\mathbf{t}(D) = D$  is obviously sufficient; however, a much simpler sufficient statistic, with a fixed dimensionality which is independent of the sample size, often exists. In fact this is known to be the case whenever the probability model belongs to the *generalized exponential family*, which includes many of the more frequently used probability models. It is easily established that if  $\mathbf{t}$  is sufficient, the posterior distribution of  $\omega$  only depends on the data  $D$  through  $\mathbf{t}(D)$ , and may be directly computed in terms of  $p(\mathbf{t} | \omega)$ , so that,  $p(\omega | D) = p(\omega | \mathbf{t}) \propto p(\mathbf{t} | \omega) p(\omega)$ .

Naturally, for fixed data and model assumptions, different priors lead to different posteriors. Indeed, Bayes' theorem may be described as a data-driven probability transformation machine which maps prior distributions (describing prior knowledge) into posterior distributions (repre-

senting combined prior and data knowledge). It is important to analyze whether or not sensible changes in the prior would induce noticeable changes in the posterior. Posterior distributions based on reference “noninformative” priors play a central role in this *sensitivity analysis* context. Investigation of the sensitivity of the posterior to changes in the prior is an important ingredient of the comprehensive analysis of the sensitivity of the final results to *all* accepted assumptions which any responsible statistical study should contain.

*Example 2. (Inference on a binomial parameter).* If the data  $D$  consist of  $n$  Bernoulli observations with parameter  $\theta$  which contain  $r$  positive trials, then  $p(D | \theta, n) = \theta^r (1 - \theta)^{n-r}$ , so that  $\mathbf{t}(D) = \{r, n\}$  is sufficient. Suppose that prior knowledge about  $\theta$  is described by a Beta distribution  $\text{Be}(\theta | \alpha, \beta)$ , so that  $p(\theta | \alpha, \beta) \propto \theta^{\alpha-1} (1 - \theta)^{\beta-1}$ . Using Bayes’ theorem, the posterior density of  $\theta$  is  $p(\theta | r, n, \alpha, \beta) \propto \theta^r (1 - \theta)^{n-r} \theta^{\alpha-1} (1 - \theta)^{\beta-1} \propto \theta^{r+\alpha-1} (1 - \theta)^{n-r+\beta-1}$ , the Beta distribution  $\text{Be}(\theta | r + \alpha, n - r + \beta)$ .

Suppose, for example, that in the light of precedent surveys, available information on the proportion  $\theta$  of citizens who would vote for a particular political measure in a referendum is described by a Beta distribution  $\text{Be}(\theta | 50, 50)$ , so that it is judged to be equally likely that the referendum would be won or lost, and it is judged that the probability that either side wins less than 60% of the vote is 0.95.



**Figure 2.** Prior and posterior densities of the proportion  $\theta$  of citizens that would vote in favour of a referendum.

A random survey of size 1500 is then conducted, where only 720 citizens declare to be in favour of the proposed measure. Using the results above, the corresponding posterior distribution is then  $\text{Be}(\theta | 770, 830)$ . These prior and posterior densities are plotted in Figure 2; it may be appreciated that, as one would expect, the effect of the data is to drastically reduce the initial uncertainty on the value of  $\theta$  and, hence, on the referendum outcome. More precisely,  $\Pr(\theta < 0.5 | 720, 1500, H, K) = 0.933$  (shaded region in Figure 2) so that, after the information from the survey has been included, the probability that the referendum will be lost should be judged to be about 93%. ◁

The general situation where the vector of interest is not the whole parameter vector  $\omega$ , but some function  $\theta = \theta(\omega)$  of possibly lower dimension than  $\omega$ , will now be considered. Let  $D$  be some observed data, let  $\{p(D | \omega), \omega \in \Omega\}$  be a probability model assumed to describe the probability mechanism which has generated  $D$ , let  $p(\omega)$  be a probability distribution describing any available information on the value of  $\omega$ , and let  $\theta = \theta(\omega) \in \Theta$  be a function of the original parameters over whose value inferences based on the data  $D$  are required. Any valid conclusion on the value of the *vector of interest*  $\theta$  will then be contained in its posterior probability distribution  $p(\theta | D)$  which is conditional on the observed data  $D$  and will naturally also depend, although not explicitly shown in the notation, on the assumed model  $\{p(D | \omega), \omega \in \Omega\}$ , and

on the available prior information encapsulated by  $p(\boldsymbol{\omega})$ . The required posterior distribution  $p(\boldsymbol{\theta} | D)$  is found by standard use of probability calculus. Indeed, by Bayes' theorem,  $p(\boldsymbol{\omega} | D) \propto p(D | \boldsymbol{\omega}) p(\boldsymbol{\omega})$ . Moreover, let  $\boldsymbol{\lambda} = \boldsymbol{\lambda}(\boldsymbol{\omega}) \in \Lambda$  be some other function of the original parameters such that  $\boldsymbol{\psi} = \{\boldsymbol{\theta}, \boldsymbol{\lambda}\}$  is a one-to-one transformation of  $\boldsymbol{\omega}$ , and let  $J(\boldsymbol{\omega}) = (\partial\boldsymbol{\psi}/\partial\boldsymbol{\omega})$  be the corresponding Jacobian matrix. Naturally, the introduction of  $\boldsymbol{\lambda}$  is not necessary if  $\boldsymbol{\theta}(\boldsymbol{\omega})$  is a one-to-one transformation of  $\boldsymbol{\omega}$ . Using standard change-of-variable probability techniques, the posterior density of  $\boldsymbol{\psi}$  is

$$p(\boldsymbol{\psi} | D) = p(\boldsymbol{\theta}, \boldsymbol{\lambda} | D) = \left[ \frac{p(\boldsymbol{\omega} | D)}{|J(\boldsymbol{\omega})|} \right]_{\boldsymbol{\omega}=\boldsymbol{\omega}(\boldsymbol{\psi})} \quad (10)$$

and the required posterior of  $\boldsymbol{\theta}$  is the appropriate *marginal* density, obtained by integration over the *nuisance parameter*  $\boldsymbol{\lambda}$ ,

$$p(\boldsymbol{\theta} | D) = \int_{\Lambda} p(\boldsymbol{\theta}, \boldsymbol{\lambda} | D) d\boldsymbol{\lambda}. \quad (11)$$

Notice that elimination of unwanted nuisance parameters, a simple integration within the Bayesian paradigm is, however, a difficult (often polemic) problem for conventional statistics.

Sometimes, the range of possible values of  $\boldsymbol{\omega}$  is effectively restricted by contextual considerations. If  $\boldsymbol{\omega}$  is known to belong to  $\Omega_c \subset \Omega$ , the prior distribution is only positive in  $\Omega_c$  and, using Bayes' theorem, it is immediately found that the restricted posterior is

$$p(\boldsymbol{\omega} | D, \boldsymbol{\omega} \in \Omega_c) = \frac{p(\boldsymbol{\omega} | D)}{\int_{\Omega_c} p(\boldsymbol{\omega} | D)}, \quad \boldsymbol{\omega} \in \Omega_c, \quad (12)$$

and obviously vanishes if  $\boldsymbol{\omega} \notin \Omega_c$ . Thus, to incorporate a restriction on the possible values of the parameters, it suffices to *renormalize* the unrestricted posterior distribution to the set  $\Omega_c \subset \Omega$  of parameter values which satisfy the required condition. Incorporation of known constraints on the parameter values, a simple renormalization within the Bayesian paradigm, is another very difficult problem for conventional statistics.

*Example 3. (Inference on normal parameters).* Let  $D = \{x_1, \dots, x_n\}$  be a random sample from a normal distribution  $N(x | \mu, \sigma)$ . The corresponding likelihood function is immediately found to be proportional to  $\sigma^{-n} \exp[-n\{s^2 + (\bar{x} - \mu)^2\}/(2\sigma^2)]$ , with  $n\bar{x} = \sum_i x_i$ , and  $ns^2 = \sum_i (x_i - \bar{x})^2$ . It may be shown (see Section 5) that absence of initial information on the value of both  $\mu$  and  $\sigma$  may formally be described by a joint prior function which is uniform in both  $\mu$  and  $\log(\sigma)$ , that is, by the (improper) prior function  $p(\mu, \sigma) = \sigma^{-1}$ . Using Bayes' theorem, the corresponding joint posterior is

$$p(\mu, \sigma | D) \propto \sigma^{-(n+1)} \exp[-n\{s^2 + (\bar{x} - \mu)^2\}/(2\sigma^2)]. \quad (13)$$

Thus, using the Gamma integral in terms of  $\lambda = \sigma^{-2}$  to integrate out  $\sigma$ ,

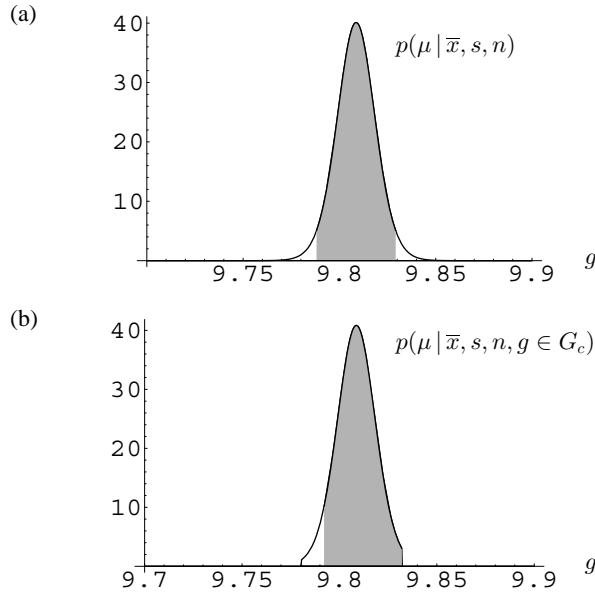
$$p(\mu | D) \propto \int_0^\infty \sigma^{-(n+1)} \exp\left[-\frac{n}{2\sigma^2}[s^2 + (\bar{x} - \mu)^2]\right] d\sigma \propto [s^2 + (\bar{x} - \mu)^2]^{-n/2}, \quad (14)$$

which is recognized as a kernel of the Student density  $\text{St}(\mu | \bar{x}, s/\sqrt{n-1}, n-1)$ . Similarly, integrating out  $\mu$ ,

$$p(\sigma | D) \propto \int_{-\infty}^\infty \sigma^{-(n+1)} \exp\left[-\frac{n}{2\sigma^2}[s^2 + (\bar{x} - \mu)^2]\right] d\mu \propto \sigma^{-n} \exp\left[-\frac{ns^2}{2\sigma^2}\right]. \quad (15)$$

Changing variables to the precision  $\lambda = \sigma^{-2}$  results in  $p(\lambda | D) \propto \lambda^{(n-3)/2} e^{ns^2\lambda/2}$ , a kernel of the Gamma density  $\text{Ga}(\lambda | (n-1)/2, ns^2/2)$ . In terms of the standard deviation  $\sigma$  this becomes  $p(\sigma | D) = p(\lambda | D) |\partial\lambda/\partial\sigma| = 2\sigma^{-3} \text{Ga}(\sigma^{-2} | (n-1)/2, ns^2/2)$ , a square-root inverted gamma density.

A frequent example of this scenario is provided by laboratory measurements made in conditions where central limit conditions apply, so that (assuming no experimental bias) those measurements may be treated as a random sample from a normal distribution centered at the quantity  $\mu$  which is being measured, and with some (unknown) standard deviation  $\sigma$ . Suppose, for example, that in an elementary physics classroom experiment to measure the gravitational field  $g$  with a pendulum, a student has obtained  $n = 20$  measurements of  $g$  yielding (in  $\text{m/sec}^2$ ) a mean  $\bar{x} = 9.8087$ , and a standard deviation  $s = 0.0428$ . Using no other information, the corresponding posterior distribution is  $p(g | D) = \text{St}(g | 9.8087, 0.0098, 19)$  represented in Figure 3(a). In particular,  $\Pr(9.788 < g < 9.829 | D) = 0.95$ , so that, with the information provided by this experiment, the gravitational field at the location of the laboratory may be expected to lie between 9.788 and 9.829 with probability 0.95.



**Figure 3.** Posterior density  $p(g | m, s, n)$  of the value  $g$  of the gravitational field, given  $n = 20$  normal measurements with mean  $m = 9.8087$  and standard deviation  $s = 0.0428$ , (a) with no additional information, and (b) with  $g$  restricted to  $G_c = \{g; 9.7803 < g < 9.8322\}$ . Shaded areas represent 95%-credible regions of  $g$ .

Formally, the posterior distribution of  $g$  should be restricted to  $g > 0$ ; however, as immediately obvious from Figure 3a, this would not have any appreciable effect, due to the fact that the likelihood function is actually concentrated on positive  $g$  values.

Suppose now that the student is further instructed to incorporate into the analysis the fact that the value of the gravitational field  $g$  at the laboratory is known to lie between 9.7803  $\text{m/sec}^2$  (average value at the Equator) and 9.8322  $\text{m/sec}^2$  (average value at the poles). The updated posterior distribution will be

$$p(g | D, g \in G_c) = \frac{\text{St}(g | m, s/\sqrt{n-1}, n)}{\int_{g \in G_c} \text{St}(g | m, s/\sqrt{n-1}, n)}, \quad g \in G_c, \quad (16)$$

represented in Figure 3(b), where  $G_c = \{g; 9.7803 < g < 9.8322\}$ . One-dimensional numerical integration may be used to verify that  $\Pr(g > 9.792 | D, g \in G_c) = 0.95$ . Moreover, if inferences about the standard deviation  $\sigma$  of the measurement procedure are also requested, the corresponding posterior distribution is found to be  $p(\sigma | D) = 2\sigma^{-3}\text{Ga}(\sigma^{-2} | 9.5, 0.0183)$ . This has a mean  $E[\sigma | D] = 0.0458$  and yields  $\Pr(0.0334 < \sigma < 0.0642 | D) = 0.95$ .  $\triangleleft$



### 3.2. Predictive Distributions

Let  $D = \{x_1, \dots, x_n\}$ ,  $x_i \in X$ , be a set of exchangeable observations, and consider now a situation where it is desired to predict the value of a future observation  $x \in X$  generated by the same random mechanism that has generated the data  $D$ . It follows from the foundations arguments discussed in Section 2 that the solution to this prediction problem is simply encapsulated by the *predictive* distribution  $p(x | D)$  describing the uncertainty on the value that  $x$  will take, given the information provided by  $D$  and any other available knowledge. Suppose that contextual information suggests the assumption that data  $D$  may be considered to be a random sample from a distribution in the family  $\{p(x | \omega), \omega \in \Omega\}$ , and let  $p(\omega)$  be a prior distribution describing available information on the value of  $\omega$ . Since  $p(x | \omega, D) = p(x | \omega)$ , it then follows from standard probability theory that  $p(x | D) = \int_{\Omega} p(x | \omega) p(\omega | D) d\omega$ , which is an average of the probability distributions of  $x$  conditional on the (unknown) value of  $\omega$ , weighted with the posterior distribution of  $\omega$  given  $D$ .

If the assumptions on the probability model are correct, the posterior predictive distribution  $p(x | D)$  will converge, as the sample size increases, to the distribution  $p(x | \omega)$  which has generated the data. Indeed, the best technique to assess the quality of the inferences about  $\omega$  encapsulated in  $p(\omega | D)$  is to check against the observed data the predictive distribution  $p(x | D)$  generated by  $p(\omega | D)$ .

*Example 4. (Prediction in a Poisson process).* Let  $D = \{r_1, \dots, r_n\}$  be a random sample from a Poisson distribution  $\text{Pn}(r | \lambda)$  with parameter  $\lambda$ , so that  $p(D | \lambda) \propto \lambda^t e^{-\lambda n}$ , where  $t = \sum r_i$ . It may be shown (see Section 5) that absence of initial information on the value of  $\lambda$  may be formally described by the (improper) prior function  $p(\lambda) = \lambda^{-1/2}$ . Using Bayes' theorem, the corresponding posterior is

$$p(\lambda | D) \propto \lambda^t e^{-\lambda n} \lambda^{-1/2} \propto \lambda^{t-1/2} e^{-\lambda n}, \quad (17)$$

the kernel of a Gamma density  $\text{Ga}(\lambda |, t + 1/2, n)$ , with mean  $(t + 1/2)/n$ . The corresponding predictive distribution is the Poisson-Gamma mixture

$$p(r | D) = \int_0^{\infty} \text{Pn}(r | \lambda) \text{Ga}(\lambda |, t + \frac{1}{2}, n) d\lambda = \frac{n^{t+1/2}}{\Gamma(t + 1/2)} \frac{1}{r!} \frac{\Gamma(r + t + 1/2)}{(1 + n)^{r+t+1/2}}. \quad (18)$$

Suppose, for example, that in a firm producing automobile restraint systems, the entire production in each of 10 consecutive months has yielded no complaint from their clients. With no additional information on the average number  $\lambda$  of complaints per month, the quality assurance department of the firm may report that the probabilities that  $r$  complaints will be received in the next month of production are given by equation (18), with  $t = 0$  and  $n = 10$ . In particular,  $p(r = 0 | D) = 0.953$ ,  $p(r = 1 | D) = 0.043$ , and  $p(r = 2 | D) = 0.003$ . Many other situations may be described with the same model. For instance, if meteorological conditions remain similar in a given area,  $p(r = 0 | D) = 0.953$  would describe the chances of no flash flood next year, given 10 years without flash floods in the area. ◁

*Example 5. (Prediction in a Normal process).* Consider now prediction of a continuous variable. Let  $D = \{x_1, \dots, x_n\}$  be a random sample from a normal distribution  $N(x | \mu, \sigma)$ . As mentioned in Example 3, absence of initial information on the values of both  $\mu$  and  $\sigma$  is formally described by the *improper* prior function  $p(\mu, \sigma) = \sigma^{-1}$ , and this leads to the joint posterior density (13). The corresponding (posterior) predictive distribution is

$$p(x | D) = \int_0^{\infty} \int_{-\infty}^{\infty} N(x | \mu, \sigma) p(\mu, \sigma | D) d\mu d\sigma = \text{St}(x | \bar{x}, s\sqrt{\frac{n+1}{n-1}}, n-1). \quad (19)$$

If  $\mu$  is known to be positive, the appropriate prior function will be the restricted function

$$p(\mu, \sigma) = \begin{cases} \sigma^{-1} & \text{if } \mu > 0 \\ 0 & \text{otherwise.} \end{cases} \quad (20)$$

However, the result in equation (19) will still hold, provided the likelihood function  $p(D | \mu, \sigma)$  is concentrated on positive  $\mu$  values. Suppose, for example, that in the firm producing automobile restraint systems, the observed breaking strengths of  $n = 10$  randomly chosen safety belt webbings have mean  $\bar{x} = 28.011$  kN and standard deviation  $s = 0.443$  kN, and that the relevant engineering specification requires breaking strengths to be larger than 26 kN. If data may truly be assumed to be a random sample from a normal distribution, the likelihood function is only appreciable for positive  $\mu$  values, and only the information provided by this small sample is to be used, then the quality engineer may claim that the probability that a safety belt randomly chosen from the same batch as the sample tested would satisfy the required specification is  $\Pr(x > 26 | D) = 0.9987$ . Besides, if production conditions remain constant, 99.87% of the safety belt webbings may be expected to have acceptable breaking strengths.  $\triangleleft$

### 3.3. Asymptotic Behaviour

The behaviour of posterior distributions when the sample size is large is now considered. This is important for, at least, two different reasons: (i) asymptotic results provide useful first-order approximations when actual samples are relatively large, and (ii) objective Bayesian methods typically depend on the asymptotic properties of the assumed model. Let  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ ,  $\mathbf{x} \in X$ , be a random sample of size  $n$  from  $\{p(\mathbf{x} | \omega), \omega \in \Omega\}$ . It may be shown that, as  $n \rightarrow \infty$ , the posterior distribution  $p(\omega | D)$  of a *discrete* parameter  $\omega$  typically converges to a degenerate distribution which gives probability one to the true value of  $\omega$ , and that the posterior distribution of a *continuous* parameter  $\omega$  typically converges to a normal distribution centered at its *maximum likelihood estimate*  $\hat{\omega}$  (MLE for short), with a variance matrix which decreases with  $n$  as  $1/n$ .

Consider first the situation where  $\Omega = \{\omega_1, \omega_2, \dots\}$  consists of a *countable* (possibly infinite) set of values, such that the probability model which corresponds to the true parameter value  $\omega_t$  is *distinguishable* from the others in the sense that the logarithmic divergence  $\delta\{p(\mathbf{x} | \omega_i) | p(\mathbf{x} | \omega_t)\}$  of each of the  $p(\mathbf{x} | \omega_i)$  from  $p(\mathbf{x} | \omega_t)$  is strictly positive. Taking logarithms in Bayes' theorem, defining  $z_j = \log[p(\mathbf{x}_j | \omega_i)/p(\mathbf{x}_j | \omega_t)]$ ,  $j = 1, \dots, n$ , and using the strong law of large numbers on the  $n$  conditionally independent and identically distributed random quantities  $z_1, \dots, z_n$ , it may be shown that

$$\lim_{n \rightarrow \infty} p(\omega_t | \mathbf{x}_1, \dots, \mathbf{x}_n) = 1, \quad \lim_{n \rightarrow \infty} p(\omega_i | \mathbf{x}_1, \dots, \mathbf{x}_n) = 0, \quad i \neq t. \quad (21)$$

Thus, under appropriate regularity conditions, the posterior probability of the true parameter value converges to one as the sample size grows.

Consider now the situation where  $\omega$  is a  $k$ -dimensional *continuous* parameter. Expressing Bayes' theorem as  $p(\omega | \mathbf{x}_1, \dots, \mathbf{x}_n) \propto \exp\{\log[p(\omega)] + \sum_{j=1}^n \log[p(\mathbf{x}_j | \omega)]\}$ , expanding  $\sum_j \log[p(\mathbf{x}_j | \omega)]$  about its maximum (the MLE  $\hat{\omega}$ ), and assuming regularity conditions (to ensure that terms of order higher than quadratic may be ignored and that the sum of the terms from the likelihood will dominate the term from the prior) it is found that the posterior density of  $\omega$  is the approximate  $k$ -variate normal

$$p(\omega | \mathbf{x}_1, \dots, \mathbf{x}_n) \approx N_k\{\hat{\omega}, \mathbf{S}(D, \hat{\omega})\}, \quad \mathbf{S}^{-1}(D, \omega) = \left( - \sum_{l=1}^n \frac{\partial^2 \log[p(\mathbf{x}_l | \omega)]}{\partial \omega_i \partial \omega_j} \right). \quad (22)$$

A simpler, but somewhat poorer, approximation may be obtained by using the strong law of large numbers on the sums in (22) to establish that  $\mathbf{S}^{-1}(D, \hat{\boldsymbol{\omega}}) \approx n \mathbf{F}(\hat{\boldsymbol{\omega}})$ , where  $\mathbf{F}(\boldsymbol{\omega})$  is Fisher's information matrix, of general element

$$\mathbf{F}_{ij}(\boldsymbol{\omega}) = - \int_X p(\mathbf{x} | \boldsymbol{\omega}) \frac{\partial^2 \log[p(\mathbf{x} | \boldsymbol{\omega})]}{\partial \omega_i \partial \omega_j} d\mathbf{x}, \quad (23)$$

so that

$$p(\boldsymbol{\omega} | \mathbf{x}_1, \dots, \mathbf{x}_n) \approx N_k(\boldsymbol{\omega} | \hat{\boldsymbol{\omega}}, n^{-1} \mathbf{F}^{-1}(\hat{\boldsymbol{\omega}})). \quad (24)$$

Thus, under appropriate regularity conditions, the posterior probability density of the parameter vector  $\boldsymbol{\omega}$  approaches, as the sample size grows, a multivariate normal density centered at the MLE  $\hat{\boldsymbol{\omega}}$ , with a variance matrix which decreases with  $n$  as  $n^{-1}$ .

*Example 2. (Inference on a binomial parameter, continued).* Let  $D = (x_1, \dots, x_n)$  consist of  $n$  independent Bernoulli trials with parameter  $\theta$ , so that  $p(D | \theta, n) = \theta^r (1 - \theta)^{n-r}$ . This likelihood function is maximized at  $\hat{\theta} = r/n$ , and Fisher's information function is  $F(\theta) = \theta^{-1}(1 - \theta)^{-1}$ . Thus, using the results above, the posterior distribution of  $\theta$  will be the approximate normal,

$$p(\theta | r, n) \approx N(\theta | \hat{\theta}, s(\hat{\theta})/\sqrt{n}), \quad s(\theta) = \{\theta(1 - \theta)\}^{1/2} \quad (25)$$

with mean  $\hat{\theta} = r/n$  and variance  $\hat{\theta}(1 - \hat{\theta})/n$ . This will provide a reasonable approximation to the exact posterior if (i) the prior  $p(\theta)$  is relatively "flat" in the region where the likelihood function matters, and (ii) both  $r$  and  $n$  are moderately large. If, say,  $n = 1500$  and  $r = 720$ , this leads to  $p(\theta | D) \approx N(\theta | 0.480, 0.013)$ , and to  $\Pr(\theta > 0.5 | D) \approx 0.940$ , which may be compared with the exact value  $\Pr(\theta > 0.5 | D) = 0.933$  obtained from the posterior distribution which corresponds to the prior  $\text{Be}(\theta | 50, 50)$ .  $\triangleleft$

It follows from the *joint* posterior asymptotic behaviour of  $\boldsymbol{\omega}$  and from the properties of the multivariate normal distribution that, if the parameter vector is decomposed into  $\boldsymbol{\omega} = (\boldsymbol{\theta}, \boldsymbol{\lambda})$ , and Fisher's information matrix is correspondingly partitioned, so that

$$\mathbf{F}(\boldsymbol{\omega}) = \mathbf{F}(\boldsymbol{\theta}, \boldsymbol{\lambda}) = \begin{pmatrix} \mathbf{F}_{\theta\theta}(\boldsymbol{\theta}, \boldsymbol{\lambda}) & \mathbf{F}_{\theta\lambda}(\boldsymbol{\theta}, \boldsymbol{\lambda}) \\ \mathbf{F}_{\lambda\theta}(\boldsymbol{\theta}, \boldsymbol{\lambda}) & \mathbf{F}_{\lambda\lambda}(\boldsymbol{\theta}, \boldsymbol{\lambda}) \end{pmatrix} \quad (26)$$

and

$$\mathbf{S}(\boldsymbol{\theta}, \boldsymbol{\lambda}) = \mathbf{F}^{-1}(\boldsymbol{\theta}, \boldsymbol{\lambda}) = \begin{pmatrix} \mathbf{S}_{\theta\theta}(\boldsymbol{\theta}, \boldsymbol{\lambda}) & \mathbf{S}_{\theta\lambda}(\boldsymbol{\theta}, \boldsymbol{\lambda}) \\ \mathbf{S}_{\lambda\theta}(\boldsymbol{\theta}, \boldsymbol{\lambda}) & \mathbf{S}_{\lambda\lambda}(\boldsymbol{\theta}, \boldsymbol{\lambda}) \end{pmatrix}, \quad (27)$$

then the *marginal* posterior distribution of  $\boldsymbol{\theta}$  will be

$$p(\boldsymbol{\theta} | D) \approx N\{\boldsymbol{\theta} | \hat{\boldsymbol{\theta}}, n^{-1} \mathbf{S}_{\theta\theta}(\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\lambda}})\}, \quad (28)$$

while the *conditional* posterior distribution of  $\boldsymbol{\lambda}$  given  $\boldsymbol{\theta}$  will be

$$p(\boldsymbol{\lambda} | \boldsymbol{\theta}, D) \approx N\{\boldsymbol{\lambda} | \hat{\boldsymbol{\lambda}} - \mathbf{F}_{\lambda\lambda}^{-1}(\boldsymbol{\theta}, \hat{\boldsymbol{\lambda}}) \mathbf{F}_{\lambda\theta}(\boldsymbol{\theta}, \hat{\boldsymbol{\lambda}})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}), n^{-1} \mathbf{F}_{\lambda\lambda}^{-1}(\boldsymbol{\theta}, \hat{\boldsymbol{\lambda}})\}. \quad (29)$$

Notice that  $\mathbf{F}_{\lambda\lambda}^{-1} = \mathbf{S}_{\lambda\lambda}$  if (and only if)  $\mathbf{F}$  is block diagonal, *i.e.*, if (and only if)  $\boldsymbol{\theta}$  and  $\boldsymbol{\lambda}$  are asymptotically independent.

*Example 3. (Inference on normal parameters, continued).* Let  $D = (x_1, \dots, x_n)$  be a random sample from a normal distribution  $N(x | \mu, \sigma)$ . The corresponding likelihood function  $p(D | \mu, \sigma)$  is maximized at  $(\hat{\mu}, \hat{\sigma}) = (\bar{x}, s)$ , and Fisher's information matrix is diagonal, with  $F_{\mu\mu} = \sigma^{-2}$ . Hence, the posterior distribution of  $\mu$  is *approximately*  $N(\mu | \bar{x}, s/\sqrt{n})$ ; this may be compared with the *exact* result  $p(\mu | D) = \text{St}(\mu | \bar{x}, s/\sqrt{n-1}, n-1)$  obtained previously under the assumption of no prior knowledge.  $\triangleleft$

## 4. Inference Summaries

From a Bayesian viewpoint, the final outcome of a problem of inference about *any* unknown quantity is precisely the corresponding posterior distribution. Thus, given some data  $D$  and conditions  $C$ , *all* that can be said about any function  $\omega$  of the parameters which govern the model is contained in the posterior distribution  $p(\omega | D, C)$ , and *all* that can be said about some function  $y$  of future observations from the same model is contained in its posterior predictive distribution  $p(y | D, C)$ . As mentioned before, Bayesian inference may technically be described as a decision problem where the space of available actions is the class of those posterior probability distributions of the quantity of interest which are compatible with accepted assumptions.

However, to make it easier for the user to assimilate the appropriate conclusions, it is often convenient to *summarize* the information contained in the posterior distribution by (i) providing values of the quantity of interest which, in the light of the data, are likely to be “close” to its true value and by (ii) measuring the compatibility of the results with hypothetical values of the quantity of interest which might have been suggested in the context of the investigation. In this section, those Bayesian counterparts of traditional *estimation* and *hypothesis testing* problems are briefly considered.

### 4.1. Estimation

In one or two dimensions, a graph of the posterior probability density of the quantity of interest (or the probability mass function in the discrete case) immediately conveys an intuitive, “impressionist” summary of the main conclusions which may possibly be drawn on its value. Indeed, this is greatly appreciated by users, and may be quoted as an important asset of Bayesian methods. From a plot of its posterior density, the region where (given the data) a univariate quantity of interest is likely to lie is easily distinguished. For instance, all important conclusions about the value of the gravitational field in Example 3 are qualitatively available from Figure 3. However, this does not easily extend to more than two dimensions and, besides, *quantitative* conclusions (in a simpler form than that provided by the mathematical expression of the posterior distribution) are often required.

*Point Estimation.* Let  $D$  be the available data, which are assumed to have been generated by a probability model  $\{p(D | \omega), \omega \in \Omega\}$ , and let  $\theta = \theta(\omega) \in \Theta$  be the quantity of interest. A *point estimator* of  $\theta$  is some function of the data  $\tilde{\theta} = \tilde{\theta}(D)$  which could be regarded as an appropriate proxy for the actual, unknown value of  $\theta$ . Formally, to choose a point estimate for  $\theta$  is a *decision problem*, where the action space is the class  $\Theta$  of possible  $\theta$  values. From a decision-theoretic perspective, to choose a point estimate  $\tilde{\theta}$  of some quantity  $\theta$  is a *decision* to act as though  $\tilde{\theta}$  were  $\theta$ , not to assert something about the value of  $\theta$  (although desire to assert something simple may well be the reason to obtain an estimate). As prescribed by the foundations of decision theory (Section 2), to solve this decision problem it is necessary to specify a *loss function*  $L(\tilde{\theta}, \theta)$  measuring the consequences of acting *as if* the true value of the quantity of interest was  $\tilde{\theta}$ , when it is actually  $\theta$ . The expected posterior loss if  $\tilde{\theta}$  was used is

$$\bar{L}[\tilde{\theta} | D] = \int_{\Theta} L(\tilde{\theta}, \theta) p(\theta | D) d\theta, \quad (30)$$

and the corresponding *Bayes estimator*  $\theta^*$  is that function of the data,  $\theta^* = \theta^*(D)$ , which minimizes this expectation.

*Example 6. (Conventional Bayes estimators).* For any given model and data, the Bayes estimator obviously depends on the chosen loss function. The loss function is context specific, and should be chosen in terms of the anticipated uses of the estimate; however, a number of conventional loss functions have been suggested for those situations where no particular uses are envisaged. These loss functions produce estimates which may be regarded as simple descriptions of the *location* of the posterior distribution. For example, if the loss function is quadratic, so that  $L(\tilde{\theta}, \theta) = (\tilde{\theta} - \theta)^t(\tilde{\theta} - \theta)$ , then the Bayes estimator is the *posterior mean*  $\theta^* = E[\theta | D]$ , assuming that the mean exists. Similarly, if the loss function is a zero-one function, so that  $L(\tilde{\theta}, \theta) = 0$  if  $\tilde{\theta}$  belongs to a ball or radius  $\epsilon$  centered in  $\theta$  and  $L(\tilde{\theta}, \theta) = 1$  otherwise, then the Bayes estimator  $\theta^*$  tends to the *posterior mode* as the ball radius  $\epsilon$  tends to zero, assuming that a unique mode exists. If  $\theta$  is univariate and the loss function is linear, so that  $L(\tilde{\theta}, \theta) = c_1(\tilde{\theta} - \theta)$  if  $\tilde{\theta} \geq \theta$ , and  $L(\tilde{\theta}, \theta) = c_2(\theta - \tilde{\theta})$  otherwise, then the Bayes estimator is the *posterior quantile* of order  $c_2/(c_1 + c_2)$ , so that  $\Pr[\theta < \theta^*] = c_2/(c_1 + c_2)$ . In particular, if  $c_1 = c_2$ , the Bayes estimator is the *posterior median*. The results derived for linear loss functions clearly illustrate the fact that *any* possible parameter value may turn out to be the Bayes estimator: it all depends on the loss function describing the consequences of the anticipated uses of the estimate.  $\triangleleft$

*Example 7. (Intrinsic estimation).* Conventional loss functions are typically non-invariant under reparametrization, so that the Bayes estimator  $\phi^*$  of a one-to-one transformation  $\phi = \phi(\theta)$  of the original parameter  $\theta$  is not necessarily  $\phi(\theta^*)$  (the posterior median, which *is* invariant, is an interesting exception). Moreover, conventional loss functions focus on the “distance” between the estimate  $\tilde{\theta}$  and the true value  $\theta$ , rather than on the “distance” between the probability models they label. Intrinsic losses directly focus on how different the probability *model*  $p(D | \theta, \lambda)$  is from its closest approximation within the family  $\{p(D | \tilde{\theta}, \lambda_i), \lambda_i \in \Lambda\}$ , and typically produce invariant solutions. An attractive example is the *intrinsic discrepancy*,  $d(\tilde{\theta}, \theta)$  defined as the minimum logarithmic divergence between a probability model labeled by  $\theta$  and a probability model labeled by  $\tilde{\theta}$ . When there are no nuisance parameters, this is given by

$$d(\tilde{\theta}, \theta) = \min\{\delta(\tilde{\theta} | \theta), \delta(\theta | \tilde{\theta})\}, \quad \delta(\theta_i | \theta) = \int_T p(\mathbf{t} | \theta) \log \frac{p(\mathbf{t} | \theta)}{p(\mathbf{t} | \theta_i)} dt, \quad (31)$$

where  $\mathbf{t} = \mathbf{t}(D) \in T$  is *any* sufficient statistic (which may well be the whole data set  $D$ ). The definition is easily extended to problems with nuisance parameters; in this case,

$$\delta(\theta_i | \omega) = \delta(\theta_i | \theta, \lambda) = \inf_{\lambda_i \in \Lambda} \int_T p(\mathbf{t} | \theta, \lambda) \log \frac{p(\mathbf{t} | \theta, \lambda)}{p(\mathbf{t} | \theta_i, \lambda_i)} dt \quad (32)$$

measures the logarithmic divergence from  $p(\mathbf{t} | \theta, \lambda)$  of its closest approximation with  $\theta = \theta_i$ , and the loss function  $d(\tilde{\theta}, \omega) = \min\{\delta(\tilde{\theta} | \theta, \lambda), \delta(\theta | \tilde{\theta}, \lambda)\}$  now depends on the complete parameter vector  $\omega = (\theta, \lambda)$ . Although not explicitly shown in the notation, the intrinsic discrepancy function typically depends on the sample size  $n$ ; indeed, when the data consist of a random sample  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  from some model  $p(\mathbf{x} | \theta, \lambda)$ , then

$$\delta(\theta_i | \theta, \lambda) = n \inf_{\lambda_i \in \Lambda} \int_X p(\mathbf{x} | \theta, \lambda) \log \frac{p(\mathbf{x} | \theta, \lambda)}{p(\mathbf{x} | \theta_i, \lambda_i)} d\mathbf{x}, \quad (33)$$

so that the discrepancy associated with the full model is simply  $n$  times the discrepancy which corresponds to a single observation. The intrinsic discrepancy is a symmetric, non-negative loss function with a direct interpretation in information-theoretical terms as the minimum amount of information which is expected to be necessary to distinguish between the model  $p(D | \theta, \lambda)$  and its closest approximation within the class  $\{p(D | \tilde{\theta}, \lambda_i), \lambda_i \in \Lambda\}$ . Moreover, it is invariant under one-to-one reparametrization of the parameter of interest  $\theta$ , and does not depend on the

choice of the nuisance parameter  $\lambda$ . The *intrinsic estimator* is naturally obtained by minimizing the posterior expected intrinsic discrepancy

$$\bar{d}(\tilde{\theta} | D) = \int_{\Omega} d(\tilde{\theta}, \omega) p(\omega | D) d\omega. \quad (34)$$

Since the intrinsic discrepancy is invariant under reparametrization, minimizing its posterior expectation produces invariant estimators.  $\triangleleft$

*Example 2. (Inference on a binomial parameter, continued).* In estimation of a binomial proportion  $\theta$ , given data  $D = (n, r)$  and a Beta prior  $\text{Be}(\theta | \alpha, \beta)$ , the quadratic loss Bayes estimate (the corresponding posterior mean) is  $E[\theta | D] = (r + \alpha)/(n + \alpha + \beta)$ , while the quadratic loss estimate of, say, the log-odds  $\phi(\theta) = \log[\theta/(1 - \theta)]$ , is  $E[\phi | D] = \psi(r + \alpha) - \psi(n - r + \beta)$  (where  $\psi(x) = d \log[\Gamma(x)]/dx$  is the *digamma* function), which is *not* equal to  $\phi(E[\theta | D])$ . The intrinsic loss function in this problem is

$$d(\tilde{\theta}, \theta) = n \min\{\delta(\tilde{\theta} | \theta), \delta(\theta | \tilde{\theta})\}, \quad \delta(\theta_i | \theta) = \theta \log \frac{\theta}{\theta_i} + (1 - \theta) \log \frac{1 - \theta}{1 - \theta_i}, \quad (35)$$

and the corresponding intrinsic estimator  $\theta^*$  is obtained by minimizing the expected posterior loss  $\bar{d}(\tilde{\theta} | D) = \int d(\tilde{\theta}, \theta) p(\theta | D) d\theta$ . The exact value of  $\theta^*$  may be obtained by numerical minimization, but a very good approximation is given by

$$\theta^* \approx \frac{1}{2} \frac{r + \alpha}{n + \alpha + \beta} + \frac{1}{2} \frac{e^{\psi(r + \alpha)}}{e^{\psi(r + \alpha)} + e^{\psi(n - r + \beta)}}. \quad (36)$$

Since intrinsic estimation is an invariant procedure, the intrinsic estimator of the log-odds will simply be the log-odds of the intrinsic estimator of  $\theta$ . As one would expect, when  $r + \alpha$  and  $n - r + \beta$  are both large, all Bayes estimators of any well-behaved function  $\phi(\theta)$  will cluster around  $\phi(E[\theta | D])$ .  $\triangleleft$

*Interval Estimation.* To describe the inferential content of the posterior distribution of the quantity of interest  $p(\theta | D)$  it is often convenient to quote regions  $R \subset \Theta$  of given probability under  $p(\theta | D)$ . For example, the identification of regions containing 50%, 90%, 95%, or 99% of the probability under the posterior may be sufficient to convey the general quantitative messages implicit in  $p(\theta | D)$ ; indeed, this is the intuitive basis of graphical representations of univariate distributions like those provided by boxplots. Any region  $R \subset \Theta$  such that  $\int_R p(\theta | D) d\theta = q$ , so that, given data  $D$ , the true value of  $\theta$  belongs to  $R$  with probability  $q$ , is said to be a *posterior  $q$ -credible region* of  $\theta$ . Notice that this provides a direct, immediately intuitive statement about the unknown quantity of interest  $\theta$  in probability terms, in marked contrast to the circumlocutory statements provided by frequentist confidence intervals. Clearly, for any given  $q$  there are generally infinitely many credible regions. A credible region is invariant under reparametrization; thus, for any  $q$ -credible region  $R$  of  $\theta$ ,  $\phi(R)$  is a  $q$ -credible region of  $\phi = \phi(\theta)$ . Sometimes, credible regions are selected to have minimum size (length, area, volume), resulting in *highest probability density* (HPD) regions, where all points in the region have larger probability density than all points outside. However, HPD regions are *not* invariant under reparametrization: the image  $\phi(R)$  of an HPD region  $R$  will be a credible region for  $\phi$ , but will not generally be HPD; indeed, there is no compelling reason to restrict attention to HPD credible regions. Posterior quantiles are often used to derive credible regions. Thus, if  $\theta_q = \theta_q(D)$  is the 100 $q$ % posterior quantile of  $\theta$ , then  $R = \{\theta; \theta \leq \theta_q\}$  is a one-sided, typically unique  $q$ -credible region, and it is invariant under reparametrization. Indeed,  $q$ -credible regions of the form  $R = \{\theta; \theta_{q/2} \leq \theta \leq \theta_{1-q/2}\}$  are easier to compute, and are often quoted in preference to HPD regions.

*Example 3. (Inference on normal parameters, continued).* In the numerical example about the value of the gravitational field described in Figure 3a, the interval [9.788, 9.829] in the unrestricted posterior density of  $g$  is a HPD, 95%-credible region for  $g$ . Similarly, the interval [9.7803, 9.8322] in Figure 3b is also a 95%-credible region for  $g$ , but it is not HPD.  $\triangleleft$

The concept of a credible region for a function  $\theta = \theta(\omega)$  of the parameter vector is trivially extended to prediction problems. Thus, a posterior  $q$ -credible region for  $x \in X$  is a subset  $R$  of the outcome space  $X$  with posterior predictive probability  $q$ , so that  $\int_R p(x | D) dx = q$ .

## 4.2. Hypothesis Testing

The posterior distribution  $p(\theta | D)$  of the quantity of interest  $\theta$  conveys immediate intuitive information on those values of  $\theta$  which, given the assumed model, may be taken to be *compatible* with the observed data  $D$ , namely, those with a relatively high probability density. Sometimes, a *restriction*  $\theta \in \Theta_0 \subset \Theta$  of the possible values of the quantity of interest (where  $\Theta_0$  may possibly consist of a single value  $\theta_0$ ) is suggested in the course of the investigation as deserving special consideration, either because restricting  $\theta$  to  $\Theta_0$  would greatly simplify the model, or because there are additional, context specific arguments suggesting that  $\theta \in \Theta_0$ . Intuitively, the *hypothesis*  $H_0 \equiv \{\theta \in \Theta_0\}$  should be judged to be *compatible* with the observed data  $D$  if there are elements in  $\Theta_0$  with a relatively high posterior density. However, a more precise conclusion is often required and, once again, this is made possible by adopting a decision-oriented approach. Formally, testing the hypothesis  $H_0 \equiv \{\theta \in \Theta_0\}$  is a *decision problem* where the action space has only two elements, namely to accept ( $a_0$ ) or to reject ( $a_1$ ) the proposed restriction. To solve this decision problem, it is necessary to specify an appropriate loss function,  $L(a_i, \theta)$ , measuring the consequences of accepting or rejecting  $H_0$  as a function of the actual value  $\theta$  of the vector of interest. Notice that this requires the statement of an *alternative*  $a_1$  to accepting  $H_0$ ; this is only to be expected, for an action is taken not because it is good, but because it is better than anything else that has been imagined.

Given data  $D$ , the optimal action will be to reject  $H_0$  if (and only if) the expected posterior loss of accepting,  $\int_{\Theta} L(a_0, \theta) p(\theta | D) d\theta$ , is larger than the expected posterior loss of rejecting,  $\int_{\Theta} L(a_1, \theta) p(\theta | D) d\theta$ , that is, if (and only if)

$$\int_{\Theta} [L(a_0, \theta) - L(a_1, \theta)] p(\theta | D) d\theta = \int_{\Theta} \Delta L(\theta) p(\theta | D) d\theta > 0. \quad (37)$$

Therefore, only the loss difference  $\Delta L(\theta) = L(a_0, \theta) - L(a_1, \theta)$ , which measures the *advantage* of rejecting  $H_0$  as a function of  $\theta$ , has to be specified. Thus, as common sense dictates, the hypothesis  $H_0$  should be rejected whenever the expected advantage of rejecting  $H_0$  is positive.

A crucial element in the specification of the loss function is a description of what is precisely meant by rejecting  $H_0$ . By assumption  $a_0$  means to act *as if*  $H_0$  were true, *i.e.*, as if  $\theta \in \Theta_0$ , but there are at least two obvious options for the alternative action  $a_1$ . This may either mean (i) the *negation* of  $H_0$ , that is to act as if  $\theta \notin \Theta_0$  or, alternatively, it may rather mean (ii) to reject the simplification implied by  $H_0$  and to keep the unrestricted model,  $\theta \in \Theta$ , which is true by assumption. Both options have been analyzed in the literature, although it may be argued that the problems of scientific data analysis where hypothesis testing procedures are typically used are better described by the second alternative. Indeed, an established model, identified by  $H_0 \equiv \{\theta \in \Theta_0\}$ , is often embedded into a more general model,  $\{\theta \in \Theta, \Theta_0 \subset \Theta\}$ , constructed to include possibly promising departures from  $H_0$ , and it is required to verify whether presently available data  $D$  are still compatible with  $\theta \in \Theta_0$ , or whether the extension to  $\theta \in \Theta$  is really required.

*Example 8. (Conventional hypothesis testing).* Let  $p(\boldsymbol{\theta} | D)$ ,  $\boldsymbol{\theta} \in \Theta$ , be the posterior distribution of the quantity of interest, let  $a_0$  be the decision to work under the restriction  $\boldsymbol{\theta} \in \Theta_0$  and let  $a_1$  be the decision to work under the complementary restriction  $\boldsymbol{\theta} \notin \Theta_0$ . Suppose, moreover, that the loss structure has the simple, zero-one form given by  $\{L(a_0, \boldsymbol{\theta}) = 0, L(a_1, \boldsymbol{\theta}) = 1\}$  if  $\boldsymbol{\theta} \in \Theta_0$  and, similarly,  $\{L(a_0, \boldsymbol{\theta}) = 1, L(a_1, \boldsymbol{\theta}) = 0\}$  if  $\boldsymbol{\theta} \notin \Theta_0$ , so that the *advantage*  $\Delta L(\boldsymbol{\theta})$  of rejecting  $H_0$  is 1 if  $\boldsymbol{\theta} \notin \Theta_0$  and it is  $-1$  otherwise. With this loss function it is immediately found that the optimal action is to reject  $H_0$  if (and only if)  $\Pr(\boldsymbol{\theta} \notin \Theta_0 | D) > \Pr(\boldsymbol{\theta} \in \Theta_0 | D)$ . Notice that this formulation requires that  $\Pr(\boldsymbol{\theta} \in \Theta_0) > 0$ , that is, that the hypothesis  $H_0$  has a strictly positive prior probability. If  $\boldsymbol{\theta}$  is a continuous parameter and  $\Theta_0$  has zero measure (for instance if  $H_0$  consists of a single point  $\boldsymbol{\theta}_0$ ), this requires the use of a non-regular “sharp” prior concentrating a positive probability mass on  $\Theta_0$ .  $\triangleleft$

*Example 9. (Intrinsic hypothesis testing).* Again, let  $p(\boldsymbol{\theta} | D)$ ,  $\boldsymbol{\theta} \in \Theta$ , be the posterior distribution of the quantity of interest, and let  $a_0$  be the decision to work under the restriction  $\boldsymbol{\theta} \in \Theta_0$ , but let  $a_1$  now be the decision to keep the general, unrestricted model  $\boldsymbol{\theta} \in \Theta$ . In this case, the advantage  $\Delta L(\boldsymbol{\theta})$  of rejecting  $H_0$  as a function of  $\boldsymbol{\theta}$  may safely be assumed to have the form  $\Delta L(\boldsymbol{\theta}) = d(\Theta_0, \boldsymbol{\theta}) - d^*$ , for some  $d^* > 0$ , where (i)  $d(\Theta_0, \boldsymbol{\theta})$  is some measure of the discrepancy between the assumed model  $p(D | \boldsymbol{\theta})$  and its closest approximation within the class  $\{p(D | \boldsymbol{\theta}_0), \boldsymbol{\theta}_0 \in \Theta_0\}$ , such that  $d(\Theta_0, \boldsymbol{\theta}) = 0$  whenever  $\boldsymbol{\theta} \in \Theta_0$ , and (ii)  $d^*$  is a context dependent *utility constant* which measures the (necessarily positive) advantage of being able to work with the simpler model when it is true. Choices of both  $d(\Theta_0, \boldsymbol{\theta})$  and  $d^*$  which may be appropriate for general use will now be described.

By similar reasons to those supporting its use in point estimation, an attractive choice for the function  $d(\Theta_0, \boldsymbol{\theta})$  is an appropriate extension of the intrinsic discrepancy; when there are no nuisance parameters, this is given by

$$d(\Theta_0, \boldsymbol{\theta}) = \min \left\{ \inf_{\boldsymbol{\theta}_0 \in \Theta_0} \delta(\boldsymbol{\theta}_0 | \boldsymbol{\theta}), \inf_{\boldsymbol{\theta}_0 \in \Theta_0} \delta(\boldsymbol{\theta} | \boldsymbol{\theta}_0) \right\}, \quad (38)$$

where  $\delta(\boldsymbol{\theta}_0 | \boldsymbol{\theta}) = \int_T p(\mathbf{t} | \boldsymbol{\theta}) \log \{p(\mathbf{t} | \boldsymbol{\theta}) / p(\mathbf{t} | \boldsymbol{\theta}_0)\} d\mathbf{t}$ , and  $\mathbf{t} = \mathbf{t}(D) \in T$  is *any* sufficient statistic, which may well be the whole dataset  $D$ . As before, if the data  $D = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  consist of a random sample from  $p(\mathbf{x} | \boldsymbol{\theta})$ , then

$$\delta(\boldsymbol{\theta}_0 | \boldsymbol{\theta}) = n \int_X p(\mathbf{x} | \boldsymbol{\theta}) \log \frac{p(\mathbf{x} | \boldsymbol{\theta})}{p(\mathbf{x} | \boldsymbol{\theta}_0)} d\mathbf{x}. \quad (39)$$

Naturally, the loss function  $d(\Theta_0, \boldsymbol{\theta})$  reduces to the intrinsic discrepancy  $d(\boldsymbol{\theta}_0, \boldsymbol{\theta})$  of Example 6 when  $\Theta_0$  contains a single element  $\boldsymbol{\theta}_0$ . Besides, as in the case of estimation, the definition is easily extended to problems with nuisance parameters, with

$$\delta(\boldsymbol{\theta}_0 | \boldsymbol{\theta}, \boldsymbol{\lambda}) = \inf_{\boldsymbol{\lambda}_0 \in \Lambda} \int_T p(\mathbf{t} | \boldsymbol{\theta}, \boldsymbol{\lambda}) \log \frac{p(\mathbf{t} | \boldsymbol{\theta}, \boldsymbol{\lambda})}{p(\mathbf{t} | \boldsymbol{\theta}_0, \boldsymbol{\lambda}_0)} d\mathbf{t}. \quad (40)$$

The hypothesis  $H_0$  should be rejected if the posterior expected advantage of rejecting is

$$\bar{d}(\Theta_0, D) = \int_{\Theta} d(\Theta_0, \boldsymbol{\theta}) p(\boldsymbol{\theta} | D) d\boldsymbol{\theta} > d^*, \quad (41)$$

for some  $d^* > 0$ . It is easily verified that the function  $\bar{d}(\Theta_0, D)$  is nonnegative. Moreover, if  $\phi = \phi(\boldsymbol{\theta})$  is a one-to-one transformation of  $\boldsymbol{\theta}$ , then  $\bar{d}(\phi(\Theta_0), D) = \bar{d}(\Theta_0, D)$ , so that the expected intrinsic loss of rejecting  $H_0$  is invariant under reparametrization.

It may be shown that, as the sample size increases, the expected value of  $\bar{d}(\Theta_0, D)$  under sampling tends to one when  $H_0$  is true, and tends to infinity otherwise; thus  $\bar{d}(\Theta_0, D)$  may be regarded as a continuous, positive measure of how inappropriate (in loss of information units)



it would be to simplify the model by accepting  $H_0$ . In traditional language,  $\bar{d}(\Theta_0, D)$  is a *test statistic* for  $H_0$  and the hypothesis should be rejected if the value of  $\bar{d}(\Theta_0, D)$  exceeds some *critical value*  $d^*$ . In sharp contrast to conventional hypothesis testing, this critical value  $d^*$  is found to be a context specific, positive utility constant  $d^*$ , which may precisely be described as the number of *information units* which the decision maker is prepared to lose in order to be able to work with the simpler model  $H_0$ , and does not depend on the sampling properties of the probability model. The procedure may be used with standard, continuous regular priors even in *sharp* hypothesis testing, when  $\Theta_0$  is a zero-measure set (as would be the case if  $\theta$  is continuous and  $\Theta_0$  contains a single point  $\theta_0$ ). Naturally, to implement the test, the utility constant  $d^*$  which defines the rejection region must be chosen.

All measurements are based on a comparison with a standard; comparison with the “canonical” problem of testing a value  $\mu = \mu_0$  for the mean of a normal distribution with known variance (see below) makes it possible to *calibrate* this *information scale*. Values of  $\bar{d}(\Theta_0, D)$  of about 1 should be regarded as an indication of no evidence against  $H_0$ , since the expected value of  $\bar{d}(\Theta_0, D)$  under  $H_0$  is precisely one. Values of  $\bar{d}(\Theta_0, D)$  of about 2.5, and 5 should be respectively regarded as an indication of mild evidence against  $H_0$ , and significant evidence against  $H_0$  since, in the canonical normal problem, these values correspond to the observed sample mean  $\bar{x}$  respectively lying 2 or 3 posterior standard deviations from the null value  $\mu_0$ . Notice that, in sharp contrast to frequentist hypothesis testing, where it is hazily recommended to adjust the significance level for dimensionality and sample size, this provides an absolute scale (in information units) which remains valid for any sample size and any dimensionality.  $\triangleleft$

*Example 10. (Testing the value of a normal mean).* Let the data  $D = \{x_1, \dots, x_n\}$  be a random sample from a normal distribution  $N(x | \mu, \sigma)$ , where  $\sigma$  is assumed to be known, and consider the “canonical” problem of testing whether these data are or are not compatible with some specific sharp hypothesis  $H_0 \equiv \{\mu = \mu_0\}$  on the value of the mean.

The conventional approach to this problem requires a non-regular prior which places a probability mass, say  $p_0$ , on the value  $\mu_0$  to be tested, with the remaining  $1 - p_0$  probability continuously distributed over  $\mathfrak{R}$ . If this prior is chosen to be  $p(\mu | \mu \neq \mu_0) = N(\mu | \mu_0, \sigma_0)$ , Bayes theorem may be used to obtain the corresponding posterior probability,

$$\Pr[\mu_0 | D, \lambda] = \frac{B_{01}(D, \lambda) p_0}{(1 - p_0) + p_0 B_{01}(D, \lambda)}, \quad (42)$$

$$B_{01}(D, \lambda) = \left(1 + \frac{n}{\lambda}\right)^{1/2} \exp\left[-\frac{1}{2} \frac{n}{n + \lambda} z^2\right], \quad (43)$$

where  $z = (\bar{x} - \mu_0)/(\sigma/\sqrt{n})$  measures, in standard deviations, the distance between  $\bar{x}$  and  $\mu_0$  and  $\lambda = \sigma^2/\sigma_0^2$  is the ratio of model to prior variance. The function  $B_{01}(D, \lambda)$ , a ratio of (integrated) likelihood functions, is called the *Bayes factor* in favour of  $H_0$ . With a conventional zero-one loss function,  $H_0$  should be rejected if  $\Pr[\mu_0 | D, \lambda] < 1/2$ . The choices  $p_0 = 1/2$  and  $\lambda = 1$  or  $\lambda = 1/2$ , describing particular forms of *sharp* prior knowledge, have been suggested in the literature for routine use. The conventional approach to sharp hypothesis testing deals with situations of *concentrated* prior probability; it *assumes* important prior knowledge about the value of  $\mu$  and, hence, should *not* be used unless this is an appropriate assumption. Moreover, as pointed out in the 1950’s by Bartlett, the resulting posterior probability is extremely sensitive to the specific prior specification. In most applications,  $H_0$  is really a hazily defined small region rather than a point. For moderate sample sizes, the posterior probability  $\Pr[\mu_0 | D, \lambda]$  is an *approximation* to the posterior probability  $\Pr[\mu_0 - \epsilon < \mu < \mu_0 + \epsilon | D, \lambda]$  for some small interval around  $\mu_0$  which would have been obtained from a regular, continuous prior heavily

concentrated around  $\mu_0$ ; however, this approximation *always* breaks down for sufficiently large sample sizes. One consequence (which is immediately apparent from the last two equations) is that for any *fixed* value of the pertinent statistic  $z$ , the posterior probability of the null,  $\Pr[\mu_0 | D, \lambda]$ , tends to one as  $n \rightarrow \infty$ . Far from being specific to this example, this unappealing behaviour of posterior probabilities based on sharp, non-regular priors (discovered by Lindley in the 1950's, and generally known as *Lindley's paradox*) is *always* present in the conventional Bayesian approach to *sharp* hypothesis testing.

The intrinsic approach may be used without assuming any sharp prior knowledge. The intrinsic discrepancy is  $d(\mu_0, \mu) = n(\mu - \mu_0)^2 / (2\sigma^2)$ , a simple transformation of the standardized distance between  $\mu$  and  $\mu_0$ . As later explained (Section 5), absence of initial information about the value of  $\mu$  may formally be described in this problem by the (improper) uniform prior function  $p(\mu) = 1$ ; Bayes' theorem may then be used to obtain the corresponding (proper) posterior distribution,  $p(\mu | D) = N(\mu | \bar{x}, \sigma / \sqrt{n})$ . The expected value of  $d(\mu_0, \mu)$  with respect to this posterior is  $\bar{d}(\mu_0, D) = (1 + z^2) / 2$ , where  $z = (\bar{x} - \mu_0) / (\sigma / \sqrt{n})$  is the standardized distance between  $\bar{x}$  and  $\mu_0$ . As foretold by the general theory, the expected value of  $\bar{d}(\mu_0, D)$  under repeated sampling is one if  $\mu = \mu_0$ , and increases linearly with  $n$  if  $\mu \neq \mu_0$ . Moreover, in this canonical example, to reject  $H_0$  whenever  $|z| > 2$  or  $|z| > 3$ , that is whenever  $\mu_0$  is 2 or 3 posterior standard deviations away from  $\bar{x}$ , respectively corresponds to rejecting  $H_0$  whenever  $\bar{d}(\mu_0, D)$  is larger than 2.5, or larger than 5. But the information scale is independent of the problem, so that rejecting the null whenever its expected discrepancy from the true model is larger than  $d^* = 5$  units of information is a *general* rule (and one which corresponds to the conventional '3 $\sigma$ ' rule in the canonical normal case).

If  $\sigma$  is unknown, the intrinsic discrepancy becomes

$$d(\mu_0, \mu, \sigma) = \frac{n}{2} \log \left[ 1 + \left( \frac{\mu - \mu_0}{\sigma} \right)^2 \right]. \quad (44)$$

Moreover, as mentioned before, absence of initial information about both  $\mu$  and  $\sigma$  may be described by the (improper) prior function  $p(\mu, \sigma) = \sigma^{-1}$ . The intrinsic test statistic  $\bar{d}(\mu_0, D)$  is found as the expected value of  $d(\mu_0, \mu, \sigma)$  under the corresponding joint posterior distribution; this may be exactly expressed in terms of hypergeometric functions, and is approximated by

$$\bar{d}(\mu_0, D) \approx \frac{1}{2} + \frac{n}{2} \log \left( 1 + \frac{t^2}{n} \right), \quad (45)$$

where  $t$  is the traditional statistic  $t = \sqrt{n-1}(\bar{x} - \mu_0) / s$ ,  $ns^2 = \sum_j (x_j - \bar{x})^2$ . For instance, for samples sizes 5, 30 and 1000, and using the utility constant  $d^* = 5$ , the hypothesis  $H_0$  would be rejected whenever  $|t|$  is respectively larger than 5.025, 3.240, and 3.007.  $\triangleleft$

## 5. Reference Analysis

Under the Bayesian paradigm, the outcome of any inference problem (the posterior distribution of the quantity of interest) combines the information provided by the data with relevant available prior information. In many situations, however, either the available prior information on the quantity of interest is too vague to warrant the effort required to have it formalized in the form of a probability distribution, or it is too subjective to be useful in scientific communication or public decision making. It is therefore important to be able to identify the mathematical form of a "noninformative" prior, a prior that would have a minimal effect, relative to the data, on the posterior inference. More formally, suppose that the probability mechanism which has generated the available data  $D$  is assumed to be  $p(D | \omega)$ , for some  $\omega \in \Omega$ , and that the quantity of interest is some real-valued function  $\theta = \theta(\omega)$  of the model parameter  $\omega$ . Without loss

of generality, it may be assumed that the probability model is of the form  $p(D | \theta, \boldsymbol{\lambda})$ ,  $\theta \in \Theta$ ,  $\boldsymbol{\lambda} \in \Lambda$ , where  $\boldsymbol{\lambda}$  is some appropriately chosen nuisance parameter vector. As described in Section 3, to obtain the required posterior distribution of the quantity of interest  $p(\theta | D)$  it is necessary to specify a *joint* prior  $p(\theta, \boldsymbol{\lambda})$ . It is now required to identify the form of that joint prior  $\pi_\theta(\theta, \boldsymbol{\lambda})$ , the  $\theta$ -reference prior, which would have a *minimal effect* on the corresponding posterior distribution of  $\theta$ ,

$$\pi(\theta | D) \propto \int_{\Lambda} p(D | \theta, \boldsymbol{\lambda}) \pi_\theta(\theta, \boldsymbol{\lambda}) d\boldsymbol{\lambda}, \quad (46)$$

a prior which, to use a conventional expression, “would let the data speak for themselves” about the likely value of  $\theta$ . Properly defined, reference *posterior* distributions have an important role to play in scientific communication, for they provide the answer to a central question in the sciences: conditional on the assumed model  $p(D | \theta, \boldsymbol{\lambda})$ , and on any further assumptions of the value of  $\theta$  on which there might be universal agreement, the reference posterior  $\pi(\theta | D)$  should specify what *could* be said about  $\theta$  if the only available information about  $\theta$  were some well-documented data  $D$ .

Much work has been done to formulate “reference” priors which would make the idea described above mathematically precise. This section concentrates on an information-theoretical based approach to derive reference distributions which may be argued to provide the most advanced general procedure available. In the formulation described below, far from ignoring prior knowledge, the reference posterior exploits certain well-defined features of a *possible* prior, namely those describing a situation where relevant knowledge about the quantity of interest (beyond that universally accepted) may be held to be negligible compared to the information about that quantity which repeated experimentation (from a particular data generating mechanism) might possibly provide. Reference analysis is appropriate in contexts where the set of inferences which could be drawn in this *possible* situation is considered to be pertinent.

Any statistical analysis contains a fair number of subjective elements; these include (among others) the data selected, the model assumptions, and the choice of the quantities of interest. Reference analysis may be argued to provide an “objective” Bayesian solution to statistical inference problems in precisely the same sense that conventional statistical methods claim to be “objective”: in that the solutions only depend on model assumptions and observed data. The whole topic of objective Bayesian methods is, however, subject to polemic; interested readers will find in the bibliography some pointers to the relevant literature.

## 5.1. Reference Distributions

*One parameter.* Consider the experiment which consists of the observation of data  $D$ , generated by a random mechanism  $p(D | \theta)$  which only depends on a real-valued parameter  $\theta \in \Theta$ , and let  $\mathbf{t} = \mathbf{t}(D) \in T$  be *any* sufficient statistic (which may well be the complete data set  $D$ ). In Shannon’s general information theory, the *amount of information*  $I^\theta\{T, p(\theta)\}$  which may be expected to be provided by  $D$ , or (equivalently) by  $\mathbf{t}(D)$ , about the value of  $\theta$  is defined by

$$I^\theta\{T, p(\theta)\} = \int_T \int_{\Theta} p(\mathbf{t}, \theta) \log \frac{p(\mathbf{t}, \theta)}{p(\mathbf{t})p(\theta)} d\theta d\mathbf{t} = E_{\mathbf{t}} \left[ \int_{\Theta} p(\theta | \mathbf{t}) \log \frac{p(\theta | \mathbf{t})}{p(\theta)} d\theta \right] \quad (47)$$

the expected logarithmic divergence of the prior from the posterior. This is naturally a *functional* of the prior  $p(\theta)$ : the larger the prior information, the smaller the information which the data may be expected to provide. The functional  $I^\theta\{T, p(\theta)\}$  is concave, non-negative, and invariant under one-to-one transformations of  $\theta$ . Consider now the amount of information  $I^\theta\{T^k, p(\theta)\}$  about  $\theta$  which may be expected from the experiment which consists of  $k$  conditionally independent replications  $\{\mathbf{t}_1, \dots, \mathbf{t}_k\}$  of the original experiment. As  $k \rightarrow \infty$ , such an

experiment would provide any *missing information* about  $\theta$  which could possibly be obtained within this framework; thus, as  $k \rightarrow \infty$ , the functional  $I^\theta\{T^k, p(\theta)\}$  will approach the missing information about  $\theta$  associated with the prior  $p(\theta)$ . Intuitively, a  $\theta$ -“noninformative” prior is one which *maximizes the missing information* about  $\theta$ . Formally, if  $\pi_k(\theta)$  denotes the prior density which maximizes  $I^\theta\{T^k, p(\theta)\}$  in the class  $\mathcal{P}$  of strictly positive prior distributions which are compatible with accepted assumptions on the value of  $\theta$  (which may well be the class of *all* strictly positive proper priors) then the  $\theta$ -reference prior  $\pi(\theta)$  is the limit as  $k \rightarrow \infty$  (in a sense to be made precise) of the sequence of priors  $\{\pi_k(\theta), k = 1, 2, \dots\}$ .

Notice that this limiting procedure is *not* some kind of asymptotic approximation, but an essential element of the *definition* of a reference prior. In particular, this definition implies that reference distributions only depend on the *asymptotic* behaviour of the assumed probability model, a feature which greatly simplifies their actual derivation.

*Example 11. (Maximum entropy).* If  $\theta$  may only take a *finite* number of values, so that the parameter space is  $\Theta = \{\theta_1, \dots, \theta_m\}$  and  $p(\theta) = \{p_1, \dots, p_m\}$ , with  $p_i = \Pr(\theta = \theta_i)$ , then the missing information associated to  $\{p_1, \dots, p_m\}$  may be shown to be

$$\lim_{k \rightarrow \infty} I^\theta\{T^k, p(\theta)\} = H(p_1, \dots, p_m) = - \sum_{i=1}^m p_i \log(p_i), \quad (48)$$

that is, the *entropy* of the prior distribution  $\{p_1, \dots, p_m\}$ .

Thus, in the finite case, the reference prior is that with *maximum entropy* in the class  $\mathcal{P}$  of priors compatible with accepted assumptions. Consequently, the reference prior algorithm contains “maximum entropy” priors as the particular case which obtains when the parameter space is *finite*, the *only* case where the original concept of entropy (in statistical mechanics, as a measure of uncertainty) is unambiguous and well-behaved. If, in particular,  $\mathcal{P}$  contains *all* priors over  $\{\theta_1, \dots, \theta_m\}$ , then the reference prior is the uniform prior,  $\pi(\theta) = \{1/m, \dots, 1/m\}$ .  $\triangleleft$

Formally, the *reference prior function*  $\pi(\theta)$  of a univariate parameter  $\theta$  is defined to be the limit of the sequence of the proper priors  $\pi_k(\theta)$  which maximize  $I^\theta\{T^k, p(\theta)\}$  in the precise sense that, for any value of the sufficient statistic  $\mathbf{t} = \mathbf{t}(D)$ , the *reference posterior*, the pointwise limit  $\pi(\theta | \mathbf{t})$  of the corresponding sequence of posteriors  $\{\pi_k(\theta | \mathbf{t})\}$ , may be obtained from  $\pi(\theta)$  by formal use of Bayes theorem, so that  $\pi(\theta | \mathbf{t}) \propto p(\mathbf{t} | \theta) \pi(\theta)$ .

Reference prior *functions* are often simply called reference priors, even though they are usually *not* probability distributions. They should *not* be considered as expressions of belief, but technical devices to obtain (proper) posterior distributions which are a limiting form of the posteriors which could have been obtained from possible prior beliefs which were relatively uninformative with respect to the quantity of interest when compared with the information which data could provide.

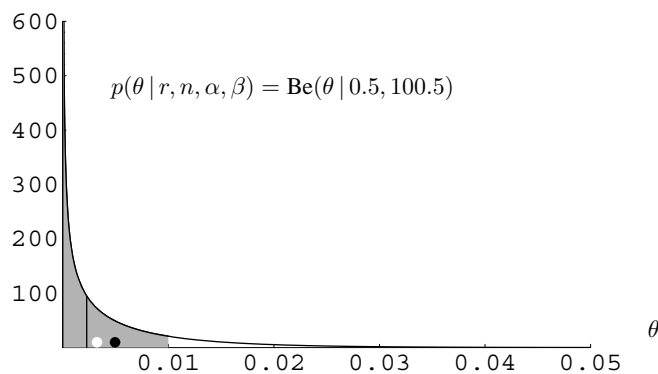
If (i) the sufficient statistic  $\mathbf{t} = \mathbf{t}(D)$  is a consistent estimator  $\tilde{\theta}$  of a continuous parameter  $\theta$ , and (ii) the class  $\mathcal{P}$  contains *all* strictly positive priors, then the reference prior may be shown to have a simple form in terms of any *asymptotic* approximation to the posterior distribution of  $\theta$ . Notice that, by construction, an *asymptotic* approximation to the posterior does *not* depend on the prior. Specifically, if the posterior density  $p(\theta | D)$  has an asymptotic approximation of the form  $p(\theta | \tilde{\theta}, n)$ , the reference prior is simply

$$\pi(\theta) \propto p(\theta | \tilde{\theta}, n) \Big|_{\tilde{\theta}=\theta} \quad (49)$$

One-parameter reference priors are shown to be *invariant* under reparametrization; thus, if  $\psi = \psi(\theta)$  is a piecewise one-to-one function of  $\theta$ , then the  $\psi$ -reference prior is simply the appropriate probability transformation of the  $\theta$ -reference prior.

*Example 12. (Jeffreys' prior).* If  $\theta$  is univariate and continuous, and the posterior distribution of  $\theta$  given  $\{x_1, \dots, x_n\}$  is asymptotically normal with standard deviation  $s(\hat{\theta})/\sqrt{n}$ , then, using (49), the reference prior function is  $\pi(\theta) \propto s(\theta)^{-1}$ . Under regularity conditions (often satisfied in practice, see Section 3.3), the posterior distribution of  $\theta$  is asymptotically normal with variance  $n^{-1} F^{-1}(\hat{\theta})$ , where  $F(\theta)$  is Fisher's information function and  $\hat{\theta}$  is the MLE of  $\theta$ . Hence, the reference prior function in these conditions is  $\pi(\theta) \propto F(\theta)^{1/2}$ , which is known as Jeffreys' prior. It follows that the reference prior algorithm contains Jeffreys' priors as the particular case which obtains when the probability model only depends on a single continuous univariate parameter, there are regularity conditions to guarantee asymptotic normality, and there is no additional information, so that the class of possible priors  $\mathcal{P}$  contains all strictly positive priors over  $\Theta$ . These are precisely the conditions under which there is general agreement on the use of Jeffreys' prior as a "noninformative" prior.  $\triangleleft$

*Example 2. (Inference on a binomial parameter, continued).* Let data  $D = \{x_1, \dots, x_n\}$  consist of a sequence of  $n$  independent Bernoulli trials, so that  $p(x|\theta) = \theta^x(1-\theta)^{1-x}$ ,  $x \in \{0, 1\}$ ; this is a regular, one-parameter continuous model, whose Fisher's information function is  $F(\theta) = \theta^{-1}(1-\theta)^{-1}$ . Thus, the reference prior  $\pi(\theta)$  is proportional to  $\theta^{-1/2}(1-\theta)^{-1/2}$ , so that the reference prior is the (proper) Beta distribution  $\text{Be}(\theta | 1/2, 1/2)$ . Since the reference algorithm is invariant under reparametrization, the reference prior of  $\phi(\theta) = 2 \arcsin \sqrt{\theta}$  is  $\pi(\phi) = \pi(\theta)/|\partial\phi/\partial\theta| = 1$ ; thus, the reference prior is *uniform on the variance-stabilizing transformation*  $\phi(\theta) = 2 \arcsin \sqrt{\theta}$ , a feature generally true under regularity conditions. In terms of the original parameter  $\theta$ , the corresponding reference posterior is  $\text{Be}(\theta | r + 1/2, n - r + 1/2)$ , where  $r = \sum x_j$  is the number of positive trials.



**Figure 4.** Posterior distribution of the proportion of infected people in the population, given the results of  $n = 100$  tests, none of which were positive.

Suppose, for example, that  $n = 100$  randomly selected people have been tested for an infection and that all tested negative, so that  $r = 0$ . The reference posterior distribution of the proportion  $\theta$  of people infected is then the Beta distribution  $\text{Be}(\theta | 0.5, 100.5)$ , represented in Figure 4. It may well be known that the infection was rare, leading to assume that  $\theta < \theta_0$ , for some upper bound  $\theta_0$ ; the (restricted) reference prior would then be of the form  $\pi(\theta) \propto \theta^{-1/2}(1-\theta)^{-1/2}$  if  $\theta < \theta_0$ , and zero otherwise. However, provided the likelihood is concentrated in the region  $\theta < \theta_0$ , the corresponding posterior would virtually be identical to  $\text{Be}(\theta | 0.5, 100.5)$ . Thus, just on the basis of the observed experimental results, one may claim that the proportion of infected people is surely smaller than 5% (for the reference posterior probability of the event  $\theta > 0.05$  is 0.001), that  $\theta$  is smaller than 0.01 with probability 0.844 (area of the shaded region in Figure 4), that it is equally likely to be over or below 0.23% (for the median, represented by a vertical line, is 0.0023), and that the probability that a person randomly chosen from the

population is infected is 0.005 (the posterior mean, represented in the figure by a black circle), since  $\Pr(x = 1 | r, n) = E[\theta | r, n] = 0.005$ . If a particular point estimate of  $\theta$  is required (say a number to be quoted in the summary headline) the *intrinsic* estimator suggests itself; this is found to be  $\theta^* = 0.0032$  (represented in the figure with a white circle). Notice that the traditional solution to this problem, based on the asymptotic behaviour of the MLE, here  $\hat{\theta} = r/n = 0$  for any  $n$ , makes absolutely no sense in this scenario.  $\triangleleft$

*One nuisance parameter.* The extension of the reference prior algorithm to the case of two parameters follows the usual mathematical procedure of reducing the problem to a sequential application of the established procedure for the single parameter case. Thus, if the probability model is  $p(\mathbf{t} | \theta, \lambda)$ ,  $\theta \in \Theta$ ,  $\lambda \in \Lambda$  and a  $\theta$ -reference prior  $\pi_\theta(\theta, \lambda)$  is required, the reference algorithm proceeds in two steps:

- (i) Conditional on  $\theta$ ,  $p(\mathbf{t} | \theta, \lambda)$  only depends on the nuisance parameter  $\lambda$  and, hence, the one-parameter algorithm may be used to obtain the *conditional* reference prior  $\pi(\lambda | \theta)$ .
- (ii) If  $\pi(\lambda | \theta)$  is proper, this may be used to integrate out the nuisance parameter thus obtaining the one-parameter integrated model  $p(\mathbf{t} | \theta) = \int_\Lambda p(\mathbf{t} | \theta, \lambda) \pi(\lambda | \theta) d\lambda$ , to which the one-parameter algorithm may be applied again to obtain  $\pi(\theta)$ . The  $\theta$ -reference prior is then  $\pi_\theta(\theta, \lambda) = \pi(\lambda | \theta) \pi(\theta)$ , and the required reference posterior is  $\pi(\theta | \mathbf{t}) \propto p(\mathbf{t} | \theta) \pi(\theta)$ .

If the conditional reference prior is *not* proper, then the procedure is performed within an increasing sequence  $\{\Lambda_i\}$  of subsets converging to  $\Lambda$  over which  $\pi(\lambda | \theta)$  is integrable. This makes it possible to obtain a corresponding sequence of  $\theta$ -reference posteriors  $\{\pi_i(\theta | \mathbf{t})\}$  for the quantity of interest  $\theta$ , and the required reference posterior is the corresponding pointwise limit  $\pi(\theta | \mathbf{t}) = \lim_i \pi_i(\theta | \mathbf{t})$ . A  $\theta$ -reference prior is then defined as a positive function  $\pi_\theta(\theta, \lambda)$  which may be formally used in Bayes' theorem as a prior to obtain the reference posterior, *i.e.*, such that, for any  $\mathbf{t} \in T$ ,  $\pi(\theta | \mathbf{t}) \propto \int_\Lambda p(\mathbf{t} | \theta, \lambda) \pi_\theta(\theta, \lambda) d\lambda$ . The approximating sequences should be *consistently* chosen within a given model. Thus, given a probability model  $\{p(x | \omega), \omega \in \Omega\}$  an appropriate approximating sequence  $\{\Omega_i\}$  should be chosen for the whole parameter space  $\Omega$ ; then, if the analysis is done in terms of, say,  $\psi = \{\psi_1, \psi_2\} \in \Psi(\Omega)$ , the approximating sequence should be chosen such that  $\Psi_i = \psi(\Omega_i)$ . A natural approximating sequence in location-scale problems is  $\{\mu, \log \sigma\} \in [-i, i]^2$ .

The  $\theta$ -reference prior does *not* depend on the choice of the nuisance parameter  $\lambda$ ; thus, for any  $\psi = \psi(\theta, \lambda)$  such that  $(\theta, \psi)$  is a one-to-one function of  $(\theta, \lambda)$ , the  $\theta$ -reference prior in terms of  $(\theta, \psi)$  is simply  $\pi_\theta(\theta, \psi) = \pi_\theta(\theta, \lambda) / |\partial(\theta, \psi) / \partial(\theta, \lambda)|$ , the appropriate probability transformation of the  $\theta$ -reference prior in terms of  $(\theta, \lambda)$ . Notice, however, that the reference prior *may* depend on the parameter of interest; thus, the  $\theta$ -reference prior may differ from the  $\phi$ -reference prior unless either  $\phi$  is a piecewise one-to-one transformation of  $\theta$ , or  $\phi$  is asymptotically independent from  $\theta$ . This is an expected consequence of the fact that the conditions under which the missing information about  $\theta$  is maximized are not generally the same as the conditions which maximize the missing information about some function  $\phi = \phi(\theta, \lambda)$ .

The *non-existence* of a unique “noninformative prior” which would be appropriate for any inference problem within a given model was established in the 1970's by Dawid and Stone, when they showed that this is incompatible with *consistent marginalization*. Indeed, if given the model  $p(D | \theta, \lambda)$ , the reference posterior of the quantity of interest  $\theta$ ,  $\pi(\theta | D) = \pi(\theta | \mathbf{t})$ , only depends on the data through a statistic  $\mathbf{t}$  whose sampling distribution,  $p(\mathbf{t} | \theta, \lambda) = p(\mathbf{t} | \theta)$ , only depends on  $\theta$ , one would expect the reference posterior to be of the form  $\pi(\theta | \mathbf{t}) \propto \pi(\theta) p(\mathbf{t} | \theta)$  for some prior  $\pi(\theta)$ . However, examples were found where this cannot be the case if a *unique* joint “noninformative” prior were to be used whatever the quantity of interest might be.

*Example 13. (Regular two dimensional continuous reference prior functions).* If the joint posterior distribution of  $(\theta, \lambda)$  is asymptotically normal, then the  $\theta$ -reference prior may be derived in terms of the corresponding Fisher's information matrix,  $\mathbf{F}(\theta, \lambda)$ . Indeed, if

$$\mathbf{F}(\theta, \lambda) = \begin{pmatrix} F_{\theta\theta}(\theta, \lambda) & F_{\theta\lambda}(\theta, \lambda) \\ F_{\theta\lambda}(\theta, \lambda) & F_{\lambda\lambda}(\theta, \lambda) \end{pmatrix}, \quad \text{and} \quad \mathbf{S}(\theta, \lambda) = \mathbf{F}^{-1}(\theta, \lambda), \quad (50)$$

then the  $\theta$ -reference prior is  $\pi_\theta(\theta, \lambda) = \pi(\lambda | \theta) \pi(\theta)$ , where

$$\pi(\lambda | \theta) \propto F_{\lambda\lambda}^{1/2}(\theta, \lambda), \quad \lambda \in \Lambda. \quad (51)$$

If  $\pi(\lambda | \theta)$  is proper,

$$\pi(\theta) \propto \exp \left\{ \int_{\Lambda} \pi(\lambda | \theta) \log[S_{\theta\theta}^{-1/2}(\theta, \lambda)] d\lambda \right\}, \quad \theta \in \Theta. \quad (52)$$

If  $\pi(\lambda | \theta)$  is not proper, integrations are performed on an approximating sequence  $\{\Lambda_i\}$  to obtain a sequence  $\{\pi_i(\lambda | \theta) \pi_i(\theta)\}$ , (where  $\pi_i(\lambda | \theta)$  is the proper renormalization of  $\pi(\lambda | \theta)$  to  $\Lambda_i$ ) and the  $\theta$ -reference prior  $\pi_\theta(\theta, \lambda)$  is defined as its appropriate limit. Moreover, if (i) both  $F_{\lambda\lambda}^{1/2}(\theta, \lambda)$  and  $S_{\theta\theta}^{-1/2}(\theta, \lambda)$  factorize, so that

$$S_{\theta\theta}^{-1/2}(\theta, \lambda) \propto f_\theta(\theta) g_\theta(\lambda), \quad F_{\lambda\lambda}^{1/2}(\theta, \lambda) \propto f_\lambda(\theta) g_\lambda(\lambda), \quad (53)$$

and (ii) the parameters  $\theta$  and  $\lambda$  are variation independent, so that  $\Lambda$  does not depend on  $\theta$ , then the  $\theta$ -reference prior is simply  $\pi_\theta(\theta, \lambda) = f_\theta(\theta) g_\lambda(\lambda)$ , even if the conditional reference prior  $\pi(\lambda | \theta) = \pi(\lambda) \propto g_\lambda(\lambda)$  (which will not depend on  $\theta$ ) is actually improper.  $\triangleleft$

*Example 3. (Inference on normal parameters, continued).* The information matrix which corresponds to a normal model  $N(x | \mu, \sigma)$  is

$$\mathbf{F}(\mu, \sigma) = \begin{pmatrix} \sigma^{-2} & 0 \\ 0 & 2\sigma^{-2} \end{pmatrix}, \quad \mathbf{S}(\mu, \sigma) = \mathbf{F}^{-1}(\mu, \sigma) = \begin{pmatrix} \sigma^2 & 0 \\ 0 & \frac{1}{2}\sigma^2 \end{pmatrix}; \quad (54)$$

hence  $F_{\sigma\sigma}^{1/2}(\mu, \sigma) = \sqrt{2} \sigma^{-1} = f_\sigma(\mu) g_\sigma(\sigma)$ , with  $g_\sigma(\sigma) = \sigma^{-1}$ , and thus  $\pi(\sigma | \mu) = \sigma^{-1}$ . Similarly,  $S_{\mu\mu}^{-1/2}(\mu, \sigma) = \sigma^{-1} = f_\mu(\mu) g_\mu(\sigma)$ , with  $f_\mu(\mu) = 1$ , and thus  $\pi(\mu) = 1$ . Therefore, the  $\mu$ -reference prior is  $\pi_\mu(\mu, \sigma) = \pi(\sigma | \mu) \pi(\mu) = \sigma^{-1}$ , as already anticipated. Moreover, as one would expect from the fact that  $\mathbf{F}(\mu, \sigma)$  is diagonal and also anticipated, it is similarly found that the  $\sigma$ -reference prior is  $\pi_\sigma(\mu, \sigma) = \sigma^{-1}$ , the same as  $\pi_\mu(\mu, \sigma)$ .

Suppose, however, that the quantity of interest is *not* the mean  $\mu$  or the standard deviation  $\sigma$ , but the *standardized* mean  $\phi = \mu/\sigma$ . Fisher's information matrix in terms of the parameters  $\phi$  and  $\sigma$  is  $\mathbf{F}(\phi, \sigma) = J^t \mathbf{F}(\mu, \sigma) J$ , where  $J = (\partial(\mu, \sigma)/\partial(\phi, \sigma))$  is the Jacobian of the inverse transformation; this yields

$$\mathbf{F}(\phi, \sigma) = \begin{pmatrix} 1 & \phi\sigma^{-1} \\ \phi\sigma^{-1} & \sigma^{-2}(2 + \phi^2) \end{pmatrix}, \quad \mathbf{S}(\phi, \sigma) = \begin{pmatrix} 1 + \frac{1}{2}\phi^2 & -\frac{1}{2}\phi\sigma \\ -\frac{1}{2}\phi\sigma & \frac{1}{2}\sigma^2 \end{pmatrix}. \quad (55)$$

Thus,  $S_{\phi\phi}^{-1/2}(\phi, \sigma) \propto (1 + \frac{1}{2}\phi^2)^{-1/2}$  and  $F_{\sigma\sigma}^{1/2}(\phi, \sigma) \propto \sigma^{-1}(2 + \phi^2)^{1/2}$ . Hence, using again the results in Example 13,  $\pi_\phi(\phi, \sigma) = (1 + \frac{1}{2}\phi^2)^{-1/2} \sigma^{-1}$ . In the original parametrization, this is  $\pi_\phi(\mu, \sigma) = (1 + \frac{1}{2}(\mu/\sigma)^2)^{-1/2} \sigma^{-2}$ , which is different from  $\pi_\mu(\mu, \sigma) = \pi_\sigma(\mu, \sigma)$ . The corresponding reference posterior of  $\phi$  is found to be  $\pi(\phi | x_1, \dots, x_n) \propto (1 + \frac{1}{2}\phi^2)^{-1/2} p(t | \phi)$  where  $t = (\sum x_j)/(\sum x_j^2)^{1/2}$ , a one-dimensional statistic whose sampling distribution,  $p(t | \mu, \sigma) = p(t | \phi)$ , only depends on  $\phi$ . Thus, the reference prior algorithm is seen to be consistent under marginalization.  $\triangleleft$

*Many parameters.* The reference algorithm is easily generalized to an arbitrary number of parameters. If the model is  $p(\mathbf{t} | \omega_1, \dots, \omega_m)$ , a joint reference prior

$$\pi(\theta_m | \theta_{m-1}, \dots, \theta_1) \times \dots \times \pi(\theta_2 | \theta_1) \times \pi(\theta_1) \quad (56)$$

may sequentially be obtained for each *ordered* parametrization  $\{\theta_1(\omega), \dots, \theta_m(\omega)\}$  of interest, and these are invariant under reparametrization of any of the  $\theta_i(\omega)$ 's. The choice of the ordered parametrization  $\{\theta_1, \dots, \theta_m\}$  precisely describes the particular prior required, namely that which *sequentially* maximizes the missing information about each of the  $\theta_i$ 's, conditional on  $\{\theta_1, \dots, \theta_{i-1}\}$ , for  $i = m, m-1, \dots, 1$ .

*Example 14. (Stein's paradox).* Let  $D$  be a random sample from a  $m$ -variate normal distribution with mean  $\boldsymbol{\mu} = \{\mu_1, \dots, \mu_m\}$  and unitary variance matrix. The reference prior which corresponds to any permutation of the  $\mu_i$ 's is uniform, and this prior leads indeed to appropriate reference posterior distributions for any of the  $\mu_i$ 's, namely  $\pi(\mu_i | D) = N(\mu_i | \bar{x}_i, 1/\sqrt{n})$ . Suppose, however, that the quantity of interest is  $\theta = \sum_i \mu_i^2$ , the distance of  $\boldsymbol{\mu}$  to the origin. As showed by Stein in the 1950's, the posterior distribution of  $\theta$  based on that uniform prior (or in any "flat" *proper* approximation) has very undesirable properties; this is due to the fact that a uniform (or nearly uniform) prior, although "noninformative" with respect to each of the individual  $\mu_i$ 's, is actually highly informative on the sum of their squares, introducing a severe positive bias (Stein's paradox). However, the reference prior which corresponds to a parametrization of the form  $\{\theta, \lambda_1, \dots, \lambda_{m-1}\}$  produces, for any choice of the nuisance parameters  $\lambda_i = \lambda_i(\boldsymbol{\mu})$ , the reference posterior  $\pi(\theta | D) = \pi(\theta | t) \propto \theta^{-1/2} \chi^2(nt | m, n\theta)$ , where  $t = \sum_i \bar{x}_i^2$ , and this posterior is shown to have the appropriate consistency properties.  $\triangleleft$

Far from being specific to Stein's example, the inappropriate behaviour in problems with many parameters of specific marginal posterior distributions derived from multivariate "flat" priors (proper or improper) is indeed very frequent. Hence, lazy, uncritical use of "flat" priors, rather than the relevant reference priors, is strongly discouraged.

*Limited information.* Although often used in contexts where no universally agreed prior knowledge about the quantity of interest is available, the reference algorithm may be used to specify a prior which incorporates any acceptable prior knowledge; it suffices to maximize the missing information within the class  $\mathcal{P}$  of priors which is compatible with such accepted knowledge. Indeed, by progressive incorporation of further restrictions into  $\mathcal{P}$ , the reference prior algorithm becomes a method of (prior) *probability assessment*. As described below, the problem has a fairly simple analytical solution when those restrictions take the form of known expected values. The incorporation of other type of restrictions usually involves numerical computations.

*Example 15. (Univariate restricted reference priors).* If the probability mechanism which is assumed to have generated the available data only depends on a univariate continuous parameter  $\theta \in \Theta \subset \mathfrak{R}$ , and the class  $\mathcal{P}$  of acceptable priors is a class of proper priors which satisfies some expected value restrictions, so that

$$\mathcal{P} = \left\{ p(\theta); \quad p(\theta) > 0, \quad \int_{\Theta} p(\theta) d\theta = 1, \quad \int_{\Theta} g_i(\theta) p(\theta) d\theta = \beta_i, \quad i = 1, \dots, m \right\} \quad (57)$$

then the (restricted) reference prior is

$$\pi(\theta | \mathcal{P}) \propto \pi(\theta) \exp \left[ \sum_{j=1}^m \gamma_j g_j(\theta) \right] \quad (58)$$

where  $\pi(\theta)$  is the unrestricted reference prior and the  $\gamma_i$ 's are constants (the corresponding Lagrange multipliers), to be determined by the restrictions which define  $\mathcal{P}$ . Suppose, for instance, that data are considered to be a random sample from a location model centered at  $\theta$ ,



and that it is further assumed that  $E[\theta] = \mu_0$  and that  $\text{Var}[\theta] = \sigma_0^2$ . The unrestricted reference prior for any regular location problem may be shown to be uniform. Thus, the restricted reference prior must be of the form  $\pi(\theta | \mathcal{P}) \propto \exp\{\gamma_1\theta + \gamma_2(\theta - \mu_0)^2\}$ , with  $\int_{\Theta} \theta \pi(\theta | \mathcal{P}) d\theta = \mu_0$  and  $\int_{\Theta} (\theta - \mu_0)^2 \pi(\theta | \mathcal{P}) d\theta = \sigma_0^2$ . Hence,  $\pi(\theta | \mathcal{P})$  is a *normal* distribution with the specified mean and variance. ◁

## 5.2. Frequentist Properties

Bayesian methods provide a *direct* solution to the problems typically posed in statistical inference; indeed, posterior distributions precisely state what can be said about unknown quantities of interest *given* available data and prior knowledge. In particular, unrestricted reference posterior distributions state what could be said if no prior knowledge about the quantities of interest were available.

A frequentist analysis of the behaviour of Bayesian procedures under repeated sampling may, however, be illuminating, for this provides some interesting bridges between frequentist and Bayesian inference. It is found that the frequentist properties of Bayesian reference procedures are typically excellent, and may be used to provide a form of calibration for reference posterior probabilities.

*Point Estimation.* It is generally accepted that, as the sample size increases, a “good” estimator  $\tilde{\theta}$  of  $\theta$  ought to get the correct value of  $\theta$  eventually, that is to be *consistent*. Under appropriate regularity conditions, any Bayes estimator  $\phi^*$  of any function  $\phi(\theta)$  converges in probability to  $\phi(\theta)$ , so that sequences of Bayes estimators are typically *consistent*. Indeed, it is known that if there is a consistent sequence of estimators, then Bayes estimators are consistent. The rate of convergence is often best for reference Bayes estimators.

It is also generally accepted that a “good” estimator should be *admissible*, that is, *not dominated* by any other estimator in the sense that its expected loss under sampling (conditional to  $\theta$ ) cannot be larger for all  $\theta$  values than that corresponding to another estimator. Any *proper* Bayes estimator is admissible; moreover, as established by Wald in the 1950’s, a procedure *must* be Bayesian (proper or improper) to be admissible. Most published admissibility results refer to quadratic loss functions, but they often extend to more general loss functions. Reference Bayes estimators are typically admissible with respect to intrinsic loss functions.

Notice, however, that many other apparently intuitive frequentist ideas on estimation have been proved to be potentially misleading. For example, given a sequence of  $n$  Bernoulli observations with parameter  $\theta$  resulting in  $r$  positive trials, the *best unbiased* estimate of  $\theta^2$  is found to be  $r(r-1)/\{n(n-1)\}$ , which yields  $\tilde{\theta}^2 = 0$  when  $r = 1$ ; but to estimate the probability of two positive trials as zero, when one positive trial has been observed, is not precisely sensible. In marked contrast, any Bayes reference estimator provides a reasonable answer. For example, the intrinsic estimator of  $\theta^2$  is simply  $(\theta^*)^2$ , where  $\theta^*$  is the intrinsic estimator of  $\theta$  described in Section 4.1. In particular, if  $r = 1$  and  $n = 2$  the intrinsic estimator of  $\theta^2$  is (as one surely might expect)  $(\theta^*)^2 = 1/4$ .

*Interval Estimation.* As the sample size increases, the frequentist coverage probability of a posterior  $q$ -credible region typically converges to  $q$  so that, for *large samples*, Bayesian credible intervals may (under regularity conditions) be interpreted as *approximate* frequentist confidence regions: under repeated sampling, a Bayesian  $q$ -credible region of  $\theta$  based on a large sample will cover the true value of  $\theta$  approximately  $100q\%$  of times. Detailed results are readily available for univariate problems. For instance, consider the probability model  $\{p(D | \omega), \omega \in \Omega\}$ , let  $\theta = \theta(\omega)$  be any univariate quantity of interest, and let  $t = t(D) \in T$  be any sufficient statistic.

If  $\theta_q(\mathbf{t})$  denotes the 100 $q$ % quantile of the posterior distribution of  $\theta$  which corresponds to some unspecified prior, so that

$$\Pr[\theta \leq \theta_q(\mathbf{t}) | \mathbf{t}] = \int_{\theta \leq \theta_q(\mathbf{t})} p(\theta | \mathbf{t}) d\theta = q, \quad (59)$$

then the coverage probability of the  $q$ -credible interval  $\{\theta; \theta \leq \theta_q(\mathbf{t})\}$ ,

$$\Pr[\theta_q(\mathbf{t}) \geq \theta | \boldsymbol{\omega}] = \int_{\theta_q(\mathbf{t}) \geq \theta} p(\mathbf{t} | \boldsymbol{\omega}) dt, \quad (60)$$

typically satisfies that  $\Pr[\theta_q(\mathbf{t}) \geq \theta | \boldsymbol{\omega}] = \Pr[\theta \leq \theta_q(\mathbf{t}) | \mathbf{t}] + O(n^{-1/2})$ . This *asymptotic* approximation is true for *all* (sufficiently regular) positive priors. However, the approximation is better, actually  $O(n^{-1})$ , for a particular class of priors known as (first-order) *probability matching* priors. Reference priors are typically found to be probability matching priors, so that they provide this improved asymptotic agreement. As a matter of fact, the agreement (in regular problems) is typically quite good even for relatively small samples.

*Example 16. (Product of normal means).* Consider the case where independent random samples  $\{x_1, \dots, x_n\}$  and  $\{y_1, \dots, y_m\}$  have respectively been taken from the normal densities  $N(x | \omega_1, 1)$  and  $N(y | \omega_2, 1)$ , and suppose that the quantity of interest is the product of their means,  $\phi = \omega_1 \omega_2$  (for instance, one may be interested in inferences about the area  $\phi$  of a rectangular piece of land, given measurements  $\{x_i\}$  and  $\{y_j\}$  of its sides). Notice that this is a simplified version of a very frequent problem in the sciences, where one is interested in the product of several magnitudes, all of which have been measured with error. Using the procedure described in Example 13, with the natural approximating sequence induced by  $(\omega_1, \omega_2) \in [-i, i]^2$ , the  $\phi$ -reference prior is found to be

$$\pi_\phi(\omega_1, \omega_2) \propto (n\omega_1^2 + m\omega_2^2)^{-1/2}, \quad (61)$$

very different from the uniform prior  $\pi_{\omega_1}(\omega_1, \omega_2) = \pi_{\omega_2}(\omega_1, \omega_2) = 1$  which should be used to make objective inferences about either  $\omega_1$  or  $\omega_2$ . The prior  $\pi_\phi(\omega_1, \omega_2)$  may be shown to provide approximate agreement between Bayesian credible regions and frequentist confidence intervals for  $\phi$ ; indeed, this prior was originally suggested by Stein in the 1980's precisely to obtain such approximate agreement. The same example was later used by Efron to stress the fact that, even within a fixed probability model  $\{p(D | \boldsymbol{\omega}), \boldsymbol{\omega} \in \Omega\}$ , the prior required to make objective inferences about some function of the parameters  $\phi = \phi(\boldsymbol{\omega})$  must generally depend on the function  $\phi$ . ◁

The numerical agreement between reference Bayesian credible regions and frequentist confidence intervals is actually perfect in special circumstances. Indeed, as Lindley pointed out in the 1950's, this is the case in those problems of inference which may be transformed to location-scale problems.

*Example 3. (Inference on normal parameters, continued).* Let  $D = \{x_1, \dots, x_n\}$  be a random sample from a normal distribution  $N(x | \mu, \sigma)$ . As mentioned before, the reference posterior of the quantity of interest  $\mu$  is the Student distribution  $\text{St}(\mu | \bar{x}, s/\sqrt{n-1}, n-1)$ . Thus, normalizing  $\mu$ , the *posterior* distribution of  $t(\mu) = \sqrt{n-1}(\bar{x} - \mu)/s$ , as a function of  $\mu$  given  $D$ , is the standard Student  $\text{St}(t | 0, 1, n-1)$  with  $n-1$  degrees of freedom. On the other hand, this function  $t$  is recognized to be precisely the conventional  $t$  statistic, whose *sampling distribution* is well known to *also* be standard Student with  $n-1$  degrees of freedom. It follows that, *for all sample sizes*, posterior *reference* credible intervals for  $\mu$  given the data will be *numerically identical* to frequentist confidence intervals based on the sampling distribution of  $t$ .

A similar result is obtained in inferences about the variance. Thus, the reference *posterior* distribution of  $\lambda = \sigma^{-2}$  is the Gamma distribution  $\text{Ga}(\sigma^{-2} | (n-1)/2, ns^2/2)$  and, hence, the *posterior* distribution of  $r = ns^2/\sigma^2$ , as a function of  $\sigma^2$  given  $D$ , is a (central)  $\chi^2$  with  $n-1$  degrees of freedom. But the function  $r$  is recognized to be a conventional statistic for this problem, whose *sampling distribution* is well known to *also* be  $\chi^2$  with  $n-1$  degrees of freedom. It follows that, *for all sample sizes*, posterior *reference* credible intervals for  $\sigma^2$  (or any one-to-one function of  $\sigma^2$ ) given the data will be *numerically identical* to frequentist confidence intervals based on the sampling distribution of  $r$ . ◁

## 6. A Stylized Case Study

To further illustrate the main aspects of Bayesian methods, and to provide a detailed, worked out example, a simplified version of a problem in engineering is analyzed below.

To study the reliability of a large production batch,  $n$  randomly selected items were put to an expensive, destructive test, yielding  $D = \{x_1, \dots, x_n\}$  as their observed lifetimes in hours of continuous use. Context considerations suggested that the lifetime  $x_i$  of each item could be assumed to be exponential with hazard rate  $\theta$ , so that  $p(x_i | \theta) = \text{Ex}[x_i | \theta] = \theta e^{-\theta x_i}$ ,  $\theta > 0$ , and that, given  $\theta$ , the lifetimes of the  $n$  items are independent. Quality engineers were interested in information on the actual value of the hazard rate  $\theta$ , and on prediction of the lifetime  $x$  of similar items. In particular, they were interested in the compatibility of the observed data with advertised values of the hazard rate, and on the proportion of items whose lifetime could be expected to be longer than some required industrial specification.

The statistical analysis of exponential data makes use of the exponential-gamma distribution  $\text{Eg}(x | \alpha, \beta)$ , obtained as a continuous mixture of exponentials with a gamma density,

$$\text{Eg}(x | \alpha, \beta) = \int_0^\infty \theta e^{-\theta x} \text{Ga}(\theta | \alpha, \beta) d\theta = \frac{\alpha \beta^\alpha}{(x + \beta)^{\alpha+1}} \quad x \geq 0, \quad \alpha > 0, \quad \beta > 0. \quad (62)$$

This is a monotonically decreasing density with mode at zero; if  $\alpha > 1$ , it has a mean  $\text{E}[x | \alpha, \beta] = \beta/(\alpha - 1)$ . Moreover, tail probabilities have a simple expression; indeed,

$$\text{Pr}[x > t | \alpha, \beta] = \left\{ \frac{\beta}{\beta + t} \right\}^\alpha. \quad (63)$$

*Likelihood function.* Under the accepted assumptions on the mechanism which generated the data,  $p(D | \theta) = \prod_j \theta e^{-\theta x_j} = \theta^n e^{-\theta s}$ , which only depends on  $s = \sum_j x_j$ , the sum of the observations. Thus,  $\mathbf{t} = (s, n)$  is a *sufficient* statistic for this model. The corresponding MLE estimator is  $\hat{\theta} = n/s$  and Fisher's information function is  $F(\theta) = \theta^{-2}$ . Moreover, the sampling distribution of  $s$  is the Gamma distribution  $p(s | \theta) = \text{Ga}(s | n, \theta)$ .

The actual data consisted of  $n = 25$  uncensored observed lifetimes which, in thousands of hours, yielded a sum  $s = 41.574$ , hence a mean  $\bar{x} = 1.663$ , and a MLE  $\hat{\theta} = 0.601$ . The standard deviation of the observed lifetimes was 1.286 and their range was  $[0.136, 5.591]$ , showing the large variation (from a few hundred to a few thousand hours) typically observed in exponential data.

Using the results of Section 3.3, and the form of Fisher's information function given above, the *asymptotic* posterior distribution of  $\theta$  is  $p(\theta | D) \approx \text{N}(\theta | \hat{\theta}, \hat{\theta}/\sqrt{n}) = \text{N}(\theta | 0.601, 0.120)$ . This provided a first, quick approximation to the possible values of  $\theta$  which, for instance, could be expected to belong to the interval  $0.601 \pm 1.96 * 0.120$ , or  $(0.366, 0.837)$ , with probability close to 0.95.

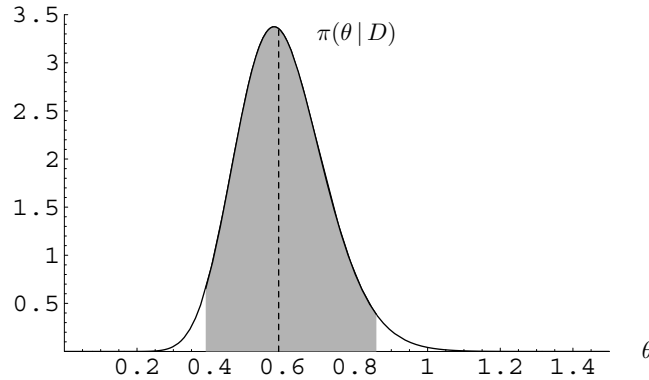
## 6.1. Objective Bayesian Analysis

The firm was to be audited on behalf of a major client. A report had to be prepared on the available information on the hazard rate  $\theta$  *exclusively* based on the *documented* data  $D$ , as if this were the *only* information available. Within a Bayesian framework, this “objective” analysis (objective in the sense of not using any information beyond that provided by the data under the assumed model) may be achieved by computing the corresponding *reference* posterior distribution.

*Reference prior and reference posteriors.* The exponential model meets all necessary regularity conditions. Thus, using the results in Example 12 and the form of Fisher’s information function mentioned above, the *reference prior function* (which in this case is also Jeffreys’ prior) is simply  $\pi(\theta) \propto F(\theta)^{1/2} = \theta^{-1}$ . Hence, using Bayes’ theorem, the reference posterior is  $\pi(\theta | D) \propto p(D|\theta) \theta^{-1} \propto \theta^{n-1} e^{-s\theta}$ , the kernel of a gamma density, so that

$$\pi(\theta | D) = \text{Ga}(\theta | n, s), \quad \theta > 0, \quad (64)$$

which has mean  $E[\theta | D] = n/s$  (which is also the MLE  $\hat{\theta}$ ), mode  $(n-1)/s$ , and standard deviation  $\sqrt{n}/s = \hat{\theta}/\sqrt{n}$ . Thus, the reference posterior of the hazard rate was found to be  $\pi(\theta | D) = \text{Ga}(\theta | 25, 41.57)$  (represented in Figure 5) with mean 0.601, mode 0.577, and standard deviation 0.120. One-dimensional numerical integration further yields  $\Pr[\theta < 0.593 | D] = 0.5$ ,  $\Pr[\theta < 0.389 | D] = 0.025$  and  $\Pr[\theta < 0.859 | D] = 0.975$ ; thus, the median is 0.593, and the interval  $(0.389, 0.859)$  is a 95% reference posterior credible region (shaded area in Figure 5). The intrinsic estimator (see below) was found to be 0.590 (dashed line in Figure 5).



**Figure 5.** Reference posterior density of the hazard rate  $\theta$ . The shaded region is a 95% credible interval. The dashed line indicates the position of the intrinsic estimator.

Under the accepted assumptions for the probability mechanism which has generated the data, the reference posterior distribution  $\pi(\theta | D) = \text{Ga}(\theta | 25, 41.57)$  contained *all* that could be said about the value of the hazard rate  $\theta$  on the exclusive basis of the observed data  $D$ . Figure 5 and the numbers quoted above respectively provided useful graphical and numerical summaries, but the fact that  $\pi(\theta | D)$  is the *complete* answer (necessary for further work on prediction or decision making) was explained to the engineers by their consultant statistician.

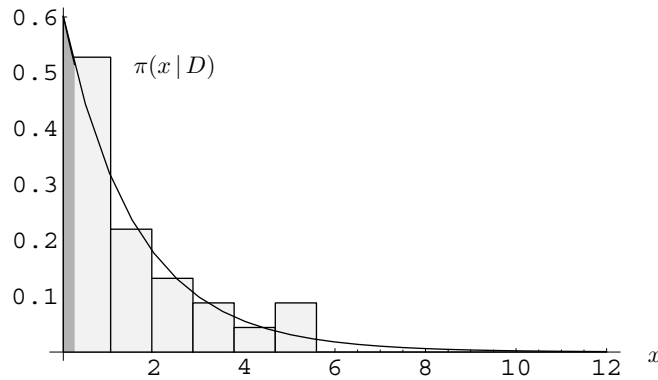
*Reference posterior predictive distribution.* The reference predictive posterior density of a future lifetime  $x$  is

$$\pi(x | D) = \int_0^{\infty} \theta e^{-\theta x} \text{Ga}(\theta | n, s) d\theta = \text{Eg}(\theta | n, s) \quad (65)$$

with mean  $s/(n - 1)$ . Thus, the posterior predictive density of the lifetime of a random item produced in similar conditions was found to be  $\pi(x | D) = \text{Eg}(x | 25, 41.57)$ , represented in Figure 6 against the background of a histogram of the observed data. The mean of this distribution is 1.732; hence, given data  $D$ , the expected lifetime of future similar items is 1.732 thousands of hours. The contract with their client specified a compensation for any item whose lifetime was smaller than 250 hours. Since

$$\Pr[x < b | D] = \int_0^b \text{Eg}(x | n, s) = 1 - \left\{ \frac{s}{s + b} \right\}^n, \quad (66)$$

the expected proportion of items with lifetime smaller than 250 hours is  $\Pr[x < 0.250 | D] = 0.139$ , the shaded area in Figure 6; thus, conditional on accepted assumptions, the engineers were advised to expect 14% of items to be nonconforming.



**Figure 6.** Reference predictive posterior density of lifetimes (in thousands of hours). The shaded region represents the probability of producing unconforming items, with lifetime smaller than 250 hours. The background is a histogram of the observed data.

*Calibration.* Consider  $t = t(\theta) = (s/n)\theta$  as a function of  $\theta$ , and its inverse transformation  $\theta = \theta(t) = (n/s)t$ . Since  $t = t(\theta)$  is a one-to-one transformation of  $\theta$ , if  $R_t$  is a  $q$ -posterior credible region for  $t$ , then  $R_\theta = \theta(R_t)$  is a  $q$ -posterior credible region for  $\theta$ . Moreover, changing variables, the reference *posterior* distribution of  $t = t(\theta)$ , as a function of  $\theta$  conditional on  $s$ , is  $\pi(t(\theta) | n, s) = \pi(\theta | n, s) / |\partial t(\theta) / \partial \theta| = \text{Ga}(t | n, n)$ , a gamma density which does not depend on  $s$ . On the other hand, the sampling distribution of the sufficient statistic  $s$  is  $p(s | n, \theta) = \text{Ga}(\theta | n, \theta)$ ; therefore, the *sampling* distribution of  $t = t(s) = (\theta/n)s$ , as a function of  $s$  conditional to  $\theta$ , is  $p(t(s) | n, \theta) = p(s | n, \theta) / |\partial t(s) / \partial s| = \text{Ga}(t | n, n)$ , which does not contain  $\theta$  and is precisely the *same* gamma density obtained before. It follows that, for *any* sample size  $n$ , all  $q$ -credible reference posterior regions of the hazard rate  $\theta$  will *also* be frequentist confidence regions of level  $q$ . Any  $q$ -credible reference posterior region has, given the data, a (rational) degree of belief  $q$  of containing the true value of  $\theta$ ; the result just obtained may be used to provide an exact calibration for this degree of belief. Indeed, for any  $\theta > 0$  and any  $q \in (0, 1)$ , the limiting proportion of  $q$ -credible reference posterior regions which would cover the true value of  $\theta$  under repeated sampling is precisely equal to  $q$ . It was therefore possible to explain to the engineers that, when reporting that the hazard rate  $\theta$  of their production was expected to be within  $(0.389, 0.859)$  with probability (rational degree of belief) 0.95, they could claim this to be a *calibrated* statement in the sense that hypothetical replications of the same *procedure* under controlled conditions, with samples simulated from *any* exponential distribution, would yield 95% of regions containing the value from which the sample was simulated.

*Estimation.* The commercial department could use any location measure of the reference posterior distribution of  $\theta$  as an intuitive estimator  $\tilde{\theta}$  of the hazard rate  $\theta$ , but if a particular value has to be chosen with, say, some legal relevance, this would pose a decision problem for which an appropriate loss function  $L(\tilde{\theta}, \theta)$  would have to be specified. Since no particular decision was envisaged, but the auditing firm nevertheless required that a particular estimator had to be quoted in the report, the attractive properties of the *intrinsic* estimator were invoked to justify its choice. The intrinsic discrepancy  $d(\theta_i, \theta_j)$  between the *models*  $\text{Ex}(x | \theta_i)$  and  $\text{Ex}(x | \theta_j)$  is

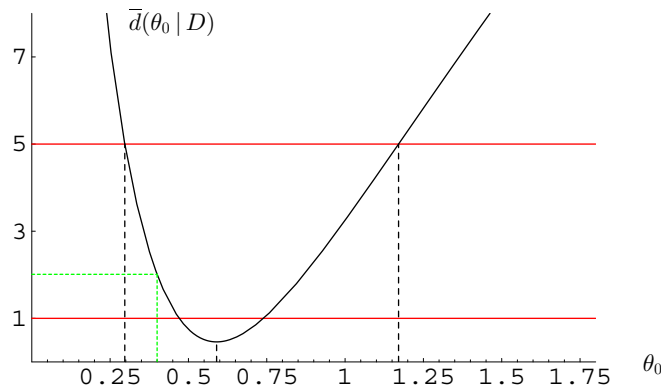
$$d(\theta_i, \theta_j) = \min\{\delta(\theta_i | \theta_j), \delta(\theta_j | \theta_i)\}, \quad \delta(\theta_i | \theta_j) = (\theta_j/\theta_i) - 1 - \log(\theta_j/\theta_i). \quad (67)$$

As expected,  $d(\theta_i, \theta_j)$  is a symmetric, non-negative concave function, which attains its minimum value zero if, and only if,  $\theta_i = \theta_j$ . The intrinsic estimator of the hazard rate is that  $\theta^*(n, s)$  which minimizes the expected reference posterior loss,

$$\bar{d}(\tilde{\theta} | n, s) = n \int_0^\infty d(\tilde{\theta}, \theta) \text{Ga}(\theta | n, s) d\theta. \quad (68)$$

To a very good approximation ( $n > 1$ ), this is given by  $\theta^*(n, s) \approx (2n - 1)/2s$ , the arithmetic average of the reference posterior mean and the reference posterior mode, quite close to the reference posterior median. With the available data, this approximation yielded  $\theta^* \approx 0.5893$ , while the exact value, found by numerical minimization was  $\theta^* = 0.5899$ . It was noticed that, since intrinsic estimation is an invariant procedure, the intrinsic estimate of any function  $\phi(\theta)$  of the hazard rate would simply be  $\phi(\theta^*)$ .

*Hypothesis Testing.* A criterion of excellence in this industrial sector described first-rate production as one with a hazard rate smaller than 0.4, yielding an expected lifetime larger than 2500 hours. The commercial department was interested in whether or not the data obtained were *compatible* with the hypothesis that the actual hazard rate of the firm's production was that small. A direct answer was provided by the corresponding reference posterior probability  $\Pr[\theta < 0.4 | D] = \int_0^{0.4} \text{Ga}(\theta | n, s) d\theta = 0.033$ , suggesting that the hazard rate of present production might possibly be around 0.4, but it is actually unlikely to be that low.



**Figure 7.** Expected reference posterior intrinsic loss for accepting  $\theta_0$  as a proxy for the true value of  $\theta$ . The minimum is reached at the intrinsic estimator  $\theta^* = 0.590$ . Values of  $\theta$  outside the interval  $(0.297, 1.170)$  would be conventionally rejected.

Under pressure to provide a quantitative measure of the compatibility of the data with the *precise* value  $\theta = \theta_0 = 0.4$ , the statistician produced the expected intrinsic discrepancy  $\bar{d}(\theta_0 | n, s)$  from accepting  $\theta_0$  as a proxy for the true value of  $\theta$  on the basis of data  $(n, s)$  by evaluating (69) at  $\tilde{\theta} = \theta_0$ . It was recalled that the expected value of  $\bar{d}(\theta_0 | D)$  under repeated sampling is precisely equal to one when  $\theta = \theta_0$ , and that a large value of  $\bar{d}(\theta_0 | D)$  indicates strong evidence

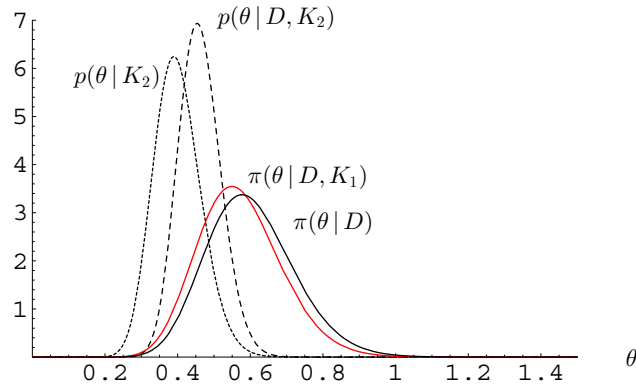
against  $\theta_0$ . Moreover, using a frequent language in engineering, the statistician explained that values of  $\bar{d}(\theta_0 | D) = d^*$  indicate, for  $d^* = 2.5, 5.0$  or  $8.5$ , a level of evidence against  $\theta = \theta_0$  comparable to the evidence against a zero mean that would be provided by a normal observation  $x$  which was, respectively, 2, 3 or 4 standard deviations from zero. As indicated in Figure 7, values of  $\theta_0$  larger than 1.170 or smaller than 0.297 would be conventionally rejected by a “ $3\sigma$ ” normal criterion. The actual value for  $\theta_0$  was found to be  $\bar{d}(0.4 | D) = 2.01$  (equivalent to  $1.73\sigma$  under normality). Thus, although there was some evidence suggesting that  $\theta$  is likely to be larger than 0.4, the precise value  $\theta = 0.4$  could not be definitely rejected on the exclusive basis of the information provided by the data  $D$ .

## 6.2. Sensitivity Analysis

Although conscious that this information could not be used in the report prepared for the client’s auditors, the firm’s management was interested in taping their engineers’ inside knowledge to gather further information on the actual lifetime of their products. This was done by exploring the consequences on the analysis of (i) introducing that information about the process which their engineers considered “beyond reasonable doubt” and (ii) introducing an “informed best guess” based on their experience with the product. The results, analyzed below and encapsulated in Figure 8, provide an analysis of the sensitivity of the final inferences on  $\theta$  to changes in the prior information.

*Limited prior information.* When questioned by their consultant statistician, the production engineers claimed to know from past experience that the average lifetime  $E[x]$  should be about 2250 hours, and that this average could not possibly be larger than 5000 or smaller than 650. Since  $E[x | \theta] = \theta^{-1}$ , those statements may directly be put in terms of conditions on the prior distribution of  $\theta$ ; indeed, working in thousands of hours, they imply  $E[\theta] = (2.25)^{-1} = 0.444$ , and that  $\theta \in \Theta_c = (0.20, 1.54)$ . To describe mathematically this knowledge  $K_1$ , the statistician used the corresponding *restricted* reference prior, that is that prior which maximizes the missing information about  $\theta$  within the class of priors which satisfy those conditions. The reference prior restricted to  $\theta \in \Theta_c$  and  $E[\theta] = \mu$  is the solution of  $\pi(\theta) \propto \theta^{-1} e^{-\lambda\theta}$ , subject to the restrictions  $\theta \in \Theta_c$  and  $\int_{\Theta_c} \theta \pi(\theta | K_1) d\theta = \mu$ . With the available data, this was numerically found to be  $\pi(\theta | K_1) \propto \theta^{-1} e^{-2.088\theta}$ ,  $\theta \in \Theta_c$ . Bayes’ theorem was then used to obtain the corresponding posterior distribution  $\pi(\theta | D, K_1) \propto p(D | \theta) \pi(\theta | K_1) \propto \theta^{24} e^{-43.69\theta}$ ,  $\theta \in \Theta_c$ , a gamma density  $\text{Ga}(\theta | 25, 43.69)$  renormalized to  $\theta \in \Theta_c$ , which is represented by a thin line in Figure 8. Comparison with the unrestricted reference posterior, described by a solid line, suggests that, compared with the information provided by the data, the additional knowledge  $K_1$  is relatively unimportant.

*Detailed prior information.* When further questioned by their consultant statistician, the production engineers guessed that the average lifetime is “surely” not larger than 3000 hours; when requested to be more precise they identified “surely” with a 0.95 subjective degree of belief. Working in thousands of hours, this implies that  $\Pr[\theta > 3^{-1}] = 0.95$ . Together with their earlier claim on the expected lifetime, implying  $E[\theta] = 0.444$ , this was sufficient to completely specify a (subjective) prior distribution  $p(\theta | K_2)$ . To obtain a tractable form for such a prior, the statistician used a simple numerical routine to fit a restricted gamma distribution to those two statements, and found this to be  $p(\theta | K_2) \propto \text{Ga}(\theta | \alpha, \beta)$ , with  $\alpha = 38.3$  and  $\beta = 86.3$ . Moreover, the statistician derived the corresponding *prior predictive* distribution  $p(x | K_2) = \text{Eg}(x | \alpha, \beta)$  and found that the elicited prior  $p(\theta)$  would imply, for instance, that  $\Pr[x > 1 | K_2] = 0.64$ ,  $\Pr[x > 3 | K_2] = 0.27$ , and  $\Pr[x > 10 | K_2] = 0.01$ , so that the implied proportion of items with a lifetime over 1, 3, and 10 thousands of hours were, respectively, 64%, 27%, and 1%. The



**Figure 8.** Probability densities of the hazard rate  $\theta$ . Subjective prior (dotted line), subjective posterior (dashed line), partially informative reference posterior (thin line) and conventional reference posterior (solid line).

engineers declared that those numbers agreed with their experience and, hence, the statistician proceeded to accept  $p(\theta) = \text{Ga}(\theta | 38.3, 86.3)$ , represented with a dotted line in Figure 8, as a reasonable description of their prior beliefs. Using Bayes' theorem, the posterior density which corresponds to a  $\text{Ga}(\theta | \alpha, \beta)$  prior is  $p(\theta | D) = p(\theta | n, s) \propto \theta^n e^{-\theta s} \theta^{\alpha-1} e^{-\beta\theta} \propto \theta^{\alpha+n-1} e^{-(\beta+s)\theta}$ , the kernel of a gamma density, so that

$$p(\theta | D) = \text{Ga}(\theta | \alpha + n, \beta + s), \quad \theta > 0. \quad (69)$$

Thus, the posterior distribution, combining the engineers' prior knowledge  $K_2$  and data  $D$  was found to be  $p(\theta | D, K_2) = \text{Ga}(\theta | 63.3, 127.8)$ , represented with a dashed line in Figure 8. It is easily appreciated from Figure 8 that the 25 observations contained in the data analyzed do not represent a dramatic increase in information over that initially claimed by the production engineers, although the posterior distribution is indeed more concentrated than the prior, and it is displaced towards the values of  $\theta$  suggested by the data. The firm's management would not be able to use this combined information in their auditing but, if they trusted their production engineers, they were advised to use  $p(\theta | D, K_2)$  to further understand their production process, or to design policies intended to improve its performance.

## 7. Discussion and Further Issues

In writing a broad article it is always hard to decide what to leave out. This article concentrates on the basic concepts of the Bayesian paradigm; methodological topics which have unwillingly been omitted include design of experiments, sample surveys, linear models and sequential methods. The interested reader is referred to the bibliography for further information. This final section briefly reviews the main arguments for the Bayesian approach, and includes pointers to further issues which have not been discussed in more detail due to space limitations.

### 7.1. Coherency

By using probability distributions to measure *all* uncertainties in the problem, the Bayesian paradigm reduces statistical inference to applied probability, thereby ensuring the coherency of the proposed solutions. There is no need to investigate, on a case by case basis, whether or not the solution to a particular problem is logically correct: a Bayesian result is only a *mathematical consequence of explicitly stated assumptions* and hence, unless a logical mistake has been committed in its derivation, it cannot be formally wrong. In marked contrast, conventional statistical methods are plagued with counterexamples. These include, among many others, negative estimators of positive quantities,  $q$ -confidence regions ( $q < 1$ ) which consist of the



whole parameter space, empty sets of “appropriate” solutions, and incompatible answers from alternative methodologies simultaneously supported by the theory.

The Bayesian approach does require, however, the specification of a (prior) probability distribution over the parameter space. The sentence “a prior distribution does not exist for this problem” is often stated to justify the use of non-Bayesian methods. However, the general representation theorem *proves the existence* of such a distribution whenever the observations are assumed to be exchangeable (and, if they are assumed to be a random sample then, *a fortiori*, they are assumed to be exchangeable). To ignore this mathematical fact, and to proceed as if a prior distribution did not exist, just because it is not easy to specify, is mathematically similar to working on a differential equation system as if no solution existed, *once it has been proved that a solution exists*, just because an explicit solution is not easily found.

## 7.2. Objectivity

It is generally accepted that any statistical analysis is subjective, in the sense that it is always conditional on accepted assumptions (on the structure of the data, on the probability model, on the outcome space) and those assumptions, although possibly well founded, are definitely *subjective* choices. It is, therefore, mandatory to make all assumptions very explicit.

Users of conventional statistical methods rarely dispute the mathematical foundations of the Bayesian approach, but claim to be able to produce “objective” answers in contrast to the possibly subjective elements involved in the choice of the prior distribution.

Bayesian methods do indeed require the choice of a prior distribution, and critics of the Bayesian approach systematically point out that in many important situations, including scientific reporting and public decision making, the results must exclusively depend on documented data which might be subject to independent scrutiny. This is of course true, but those critics choose to ignore that this particular case is covered within the Bayesian approach by the use of *reference* prior distributions which (i) are mathematically derived from the accepted probability model (and, hence, they are “objective” insofar as the choice of that model might be objective) and, (ii) by construction, they produce posterior probability distributions which, given the accepted probability model, *only* contain the information about their values which data may provide and, *optionally*, any further contextual information over which there might be universal agreement.

A related issue to that of objectivity is that of the operational meaning of reference posterior probabilities; it is found that the analysis of their behaviour under repeated sampling provides a suggestive form of calibration. Indeed,  $\Pr[\theta \in R | D] = \int_R \pi(\theta | D) d\theta$ , the reference posterior probability that  $\theta \in R$ , is *both* a measure of the conditional uncertainty (given the assumed model and the observed data  $D$ ) about the event that the unknown value of  $\theta$  belongs to  $R \subset \Theta$ , and the limiting proportion of the regions which would cover  $\theta$  under repeated sampling conditional on data “sufficiently similar” to  $D$ . Under broad conditions (to guarantee regular asymptotic behaviour), all large data sets from the same model are “sufficiently similar” among themselves in this sense and hence, given those conditions, reference posterior credible regions are *approximate* unconditional frequentist confidence regions.

The conditions for this approximate *unconditional* equivalence to hold exclude, however, important special cases, like those involving “extreme” or “relevant” observations. In very special situations, when probability models may be transformed to location-scale models, there is an exact unconditional equivalence; in those cases reference posterior credible intervals are, for any sample size, exact unconditional frequentist confidence intervals.

### 7.3. Applicability

In sharp contrast to most conventional statistical methods, which may only be exactly applied to a handful of relatively simple stylized situations, Bayesian methods are (in theory) totally general. Indeed, for a given probability model and prior distribution over its parameters, derivation of posterior distributions is a well-defined mathematical exercise. In particular, Bayesian methods do not require any particular regularity conditions on the probability model, do not depend on the existence of sufficient statistics of finite dimension, do not rely on asymptotics, and do not require the derivation of any sampling distribution, nor (a fortiori) the existence of a “pivotal” statistic whose sampling distribution is independent of the parameters.

However, when used in complex models with many parameters, Bayesian methods often require the computation of multidimensional definite integrals and, for a long time, this effectively placed practical limits on the complexity of the problems which could be handled. This has dramatically changed in recent years with the general availability of large computing power, and the parallel development of simulation-based numerical integration strategies like *importance sampling* or *Markov chain Monte Carlo* (MCMC). Those methods provide a structure within which many complex models may be analyzed using generic software. MCMC is numerical integration using Markov chains. Monte Carlo integration proceeds by drawing samples from the required distributions, and computing sample averages to approximate expectations. MCMC methods draw the required samples by running appropriately defined Markov chains for a long time; specific methods to construct those chains include the Gibbs sampler and the Metropolis algorithm, originated in the 1950’s in the literature of statistical physics. The production of improved algorithms, and the development of appropriate diagnostic tools to establish their convergence, remains a very active research area.

Actual scientific research often requires the use of models that are far too complex for conventional statistical methods to be able to handle. This article concludes with a very brief glimpse at some of them.

*Hierarchical structures.* Consider a situation where a possibly variable number  $n_i$  of observations,  $\{\mathbf{x}_{ij}, j = 1, \dots, n_i\}$ ,  $i = 1, \dots, m$ , are made on each of  $m$  internally homogeneous subsets of some population. For instance, a firm might have chosen  $m$  production lines for inspection, and  $n_i$  items might have been randomly selected among those made by production line  $i$ , so that  $\mathbf{x}_{ij}$  is the result of the measurements made on item  $j$  of production line  $i$ . As another example, animals of some species are captured to study their metabolism, and a blood sample taken before releasing them again; the procedure is repeated in the same habitat for some time, so that some of the animals are recaptured several times, and  $\mathbf{x}_{ij}$  is the result of the analysis of the  $j$ -th blood sample taken from animal  $i$ . In those situations, it is often appropriate to assume that the  $n_i$  observations on subpopulation  $i$  are exchangeable, so that they may be treated as a random sample from some model  $p(\mathbf{x} | \boldsymbol{\theta}_i)$  indexed by a parameter  $\boldsymbol{\theta}_i$  which depends on the subpopulation observed, and that the parameters which label the subpopulations may also be assumed to be exchangeable, so that  $\{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_m\}$  may be treated as a random sample from some distribution  $p(\boldsymbol{\theta} | \boldsymbol{\omega})$ . Thus, the complete *hierarchical* model which is assumed to have generated the observed data  $D = \{\mathbf{x}_{11}, \dots, \mathbf{x}_{mn_m}\}$  is of the form

$$p(D | \boldsymbol{\omega}) = \int_{\Theta^m} \left[ \prod_{j=1}^{n_i} p(\mathbf{x}_{ij} | \boldsymbol{\theta}_i) \right] \left[ \prod_{i=1}^m p(\boldsymbol{\theta}_i | \boldsymbol{\omega}) \right] \left[ \prod_{i=1}^m d\boldsymbol{\theta}_i \right]. \quad (70)$$

Hence, under the Bayesian paradigm, a family of conventional probability models, say  $p(\mathbf{x} | \boldsymbol{\theta})$ ,  $\boldsymbol{\theta} \in \Theta$ , and an appropriate “structural” prior  $p(\boldsymbol{\theta} | \boldsymbol{\omega})$ , may be naturally combined to produce a versatile, complex model  $\{p(D | \boldsymbol{\omega}), \boldsymbol{\omega} \in \Omega\}$  whose analysis is often well beyond the scope of

conventional statistics. The Bayesian solution only requires the specification a prior distribution  $p(\boldsymbol{\omega})$ , the use Bayes' theorem to obtain the corresponding posterior  $p(\boldsymbol{\omega} | D) \propto p(D | \boldsymbol{\omega}) p(\boldsymbol{\omega})$ , and the performance of the appropriate probability transformations to derive the posterior distributions of the quantities of interest (which may well be functions of  $\boldsymbol{\omega}$ , functions of the  $\theta_i$ 's, or functions of future observations). As in any other Bayesian analysis, the prior distribution  $p(\boldsymbol{\omega})$  has to describe available knowledge about  $\boldsymbol{\omega}$ ; if none is available, or if an objective analysis is required, an appropriate reference prior function  $\pi(\boldsymbol{\omega})$  may be used.

*Contextual information.* In many problems of statistical inference, objective, universally agreed, contextual information is available on the parameter values. This information is typically very difficult to handle within conventional statistics, but it is trivially incorporated into a Bayesian analysis by simply restricting the prior distribution to the class  $\{\mathcal{P}\}$  of priors which are compatible with such information. As an example, consider the frequent problem in archaeology of trying to establish the occupation period  $[\alpha, \beta]$  of a site by some past culture on the basis of the radiocarbon dating of organic samples taken from the excavation. Radiocarbon dating is not precise, so that each dating  $x_i$  is typically taken to be a normal observation from a distribution  $N(x | \mu(\theta_i), \sigma_i)$ , where  $\theta_i$  is the actual, unknown calendar date of the sample,  $\mu(\theta)$  is an internationally agreed calibration curve, and  $\sigma_i$  is a known standard error quoted by the laboratory. The actual calendar dates  $\{\theta_1, \dots, \theta_m\}$  of the samples are typically assumed to be uniformly distributed within the occupation period  $[\alpha, \beta]$ ; however, stratigraphic evidence indicates some partial orderings for, if sample  $i$  was found on top of sample  $j$  in undisturbed layers, then  $\theta_i > \theta_j$ . Thus, if  $\mathcal{C}$  denotes the class of values of  $\{\theta_1, \dots, \theta_m\}$  which satisfy those known restrictions, data may be assumed to have been generated by the hierarchical model

$$p(x_1, \dots, x_m | \alpha, \beta) = \int_{\mathcal{C}} \left[ \prod_{i=1}^m N(x_i | \mu(\theta_i), \sigma_i^2) \right] (\beta - \alpha)^{-m} d\theta_1 \dots d\theta_m. \quad (71)$$

Often, contextual information further indicates an absolute lower bound  $\alpha_0$  and an absolute upper bound  $\beta_0$  for the period investigated, so that  $\alpha_0 < \alpha < \beta < \beta_0$ . If no further documented information is available, the corresponding restricted reference prior for the quantities of interest,  $\{\alpha, \beta\}$  should be used; this is found to be  $\pi(\alpha, \beta) \propto (\beta - \alpha)^{-1}$  whenever  $\alpha_0 < \alpha < \beta < \beta_0$  and zero otherwise. The corresponding reference posterior  $\pi(\alpha, \beta | x_1, \dots, x_m) \propto p(x_1, \dots, x_m | \alpha, \beta) \pi(\alpha, \beta)$  summarizes all available information on the occupation period.

*Covariate information.* Within the last 30 years, both linear and non-linear regression models have been analyzed from a Bayesian point of view at increasing levels of sophistication. This ranges from the elementary objective Bayesian analysis of simple linear regression structures (which parallel their frequentist counterparts) to the sophisticated analysis of time series involved in dynamic forecasting which often make use of complex hierarchical structures. The field is far too large to be reviewed in this article, but the bibliography contains some relevant pointers.

*Model Criticism.* It has been stressed that *any* statistical analysis is conditional on the accepted assumptions of the probability model which is presumed to have generated the data. Recent years have shown a huge effort into the development of Bayesian procedures for *model criticism* and *model choice*. Most of these are sophisticated elaborations of the procedures described in Section 4.2 under the heading of hypothesis testing. Again, this is too large a topic to be reviewed here, but some key references are included in the bibliography.

## Acknowledgements

This work has been partially funded with Grant PB97-1403 of the DGICYT, Madrid, Spain. The author is indebted to many colleagues for their suggestions to earlier versions of this article; the detailed comments gratefully received from Dennis Lindley, Jennifer Pittman and Reinhard Viertl require, however, specific mention.

## Bibliography

- Berger, J. O. (1985). *Statistical Decision Theory and Bayesian Analysis*. Berlin: Springer. [A thorough account of Bayesian methods emphasizing its decision-theoretical aspects]
- Bernardo, J. M. (1979). Expected information as expected utility. *Ann. Statist.* **7**, 686–690. [Establishes statistical inference as a decision problem with an information-based utility function]
- Bernardo, J. M. (1997). Noninformative priors do not exist. *J. Statist. Planning and Inference* **65**, 159–189 (with discussion). [A non-technical analysis of the polemic on objective Bayesian statistics]
- Bernardo, J. M. (1999). Nested hypothesis testing: The Bayesian reference criterion. *Bayesian Statistics 6* (J. M. Bernardo, J. O. Berger, A. P. Dawid and A. F. M. Smith, eds.). Oxford: University Press, 101–130 (with discussion). [A decision-oriented approach to sharp hypothesis testing]
- Bernardo, J. M. and Ramón, J. M. (1998). An introduction to Bayesian reference analysis: inference on the ratio of multinomial parameters. *The Statistician* **47**, 1–35. [An elementary introduction to objective Bayesian analysis]
- Bernardo, J. M. and Smith, A. F. M. (1994). *Bayesian Theory*, Chichester: Wiley. [A thorough account of key concepts and theoretical results in Bayesian statistics at a graduate level, with extensive bibliography]
- Bernardo, J. M., Berger, J. O., Dawid, A. P. and Smith, A. F. M. (Eds.) (1999). *Bayesian Statistics 6*, Oxford: University Press. [The Proceedings of the 6th Valencia International Meeting on Bayesian Statistics; the Valencia meetings, held every four years, provide definite up-to-date overviews on current research within the Bayesian paradigm]
- Berry, D. A. (1996). *Statistics, a Bayesian Perspective*. Belmont, CA: Wadsworth. [A very good introduction to Bayesian statistics from a subjectivist viewpoint]
- Box, G. E. P. and Tiao, G. C. (1973). *Bayesian Inference in Statistical Analysis*. Reading, MA: Addison-Wesley. [An excellent objective Bayesian account of standard statistical problems]
- DeGroot, M. H. (1970). *Optimal Statistical Decisions*, New York: McGraw-Hill. [A thorough account of Bayesian decision theory and Bayesian inference with a rigorous treatment of foundations]
- Efron, B. (1986). Why isn't everyone a Bayesian? *Amer. Statist.* **40**, 1–11 (with discussion). [A good example of the polemic between Bayesian and non-Bayesian approaches to statistics]
- de Finetti, B. (1970). *Teoria delle Probabilità*, Turin: Einaudi. English translation as *Theory of Probability* in 1975, Chichester: Wiley. [An outstanding book on probability and statistics from a subjective viewpoint]

- Geisser, S. (1993). *Predictive Inference: an Introduction*. London: Chapman and Hall. [A comparative account of frequentist and objective Bayesian methods of prediction]
- Gelfand, A. E. and Smith, A. F. M. (1990). Sampling based approaches to calculating marginal densities. *J. Amer. Statist. Assoc.* **85**, 398-409. [An excellent primer on simulation-based techniques to numerical integration in the context of Bayesian statistics]
- Gelman, A., Carlin, J. B., Stern, H. and Rubin, D. B. (1995). *Bayesian Data Analysis*. London: Chapman and Hall. [A comprehensive treatment of Bayesian data analysis emphasizing computational tools]
- Gilks, W. R., Richardson, S. and Spiegelhalter, D. J. (1996). *Markov Chain Monte Carlo in Practice*. London: Chapman and Hall. [An excellent introduction to MCMC methods and their applications]
- Kass, R. E. and Raftery, A. E. (1995). Bayes factors. *J. Amer. Statist. Assoc.* **90**, 773–795. [A very good review of Bayes factor methods for hypothesis testing]
- Lindley, D. V. (1972). *Bayesian Statistics, a Review*. Philadelphia, PA: SIAM. [A sharp comprehensive review of the whole subject up to the 1970's, emphasizing its internal consistency]
- Lindley, D. V. (1985). *Making Decisions*. (2nd ed.) Chichester: Wiley. [The best elementary introduction to Bayesian decision analysis]
- Lindley, D. V. (1990). The 1988 Wald memorial lecture: The present position in Bayesian Statistics. *Statist. Sci.* **5**, 44-89 (with discussion). [An informative account of the Bayesian paradigm and its relationship with other attitudes to inference]
- Lindley, D. V. (2000). The philosophy of statistics. *The Statistician* **49**, 293–337 (with discussion). [A recent description of the Bayesian paradigm from a subjectivist viewpoint]
- O'Hagan, A. (1994). *Bayesian Inference* London: Edward Arnold. [A good account of Bayesian inference integrated into Kendall's Library of Statistics]
- Press, S. J. (1972). *Applied Multivariate Analysis: using Bayesian and Frequentist Methods of Inference*. Melbourne, FL: Krieger. [A comprehensive comparative account of frequentist and objective Bayesian methods of inference in multivariate problems]
- Press, S. J. and Tanur, J. M. (2001). *The Subjectivity of Scientists and the Bayesian Approach*. New York: Wiley. [An interesting description of how the preconceived beliefs of famous scientists throughout history influenced their scientific conclusions]
- West, M. and Harrison, P. J. (1989). *Bayesian Forecasting and Dynamic Models*. Berlin: Springer. [An excellent thorough account of Bayesian time series analysis]
- Zellner, A. (1971). *An Introduction to Bayesian Inference in Econometrics*. New York: Wiley. Reprinted in 1987, Melbourne, FL: Krieger. [A detailed objective Bayesian analysis of linear models]