

## **Bayesian inference in processing experimental data: principles and basic applications**

**G D'Agostini**

Dipartimento di Fisica Università 'La Sapienza' and INFN, P.le A. Moro 2, 00185 Roma, Italy

E-mail: [giulio.dagostini@roma1.infn.it](mailto:giulio.dagostini@roma1.infn.it)

Received 26 February 2003, in final form 7 July 2003

Published 11 August 2003

Online at [stacks.iop.org/RoPP/66/1383](http://stacks.iop.org/RoPP/66/1383)

### **Abstract**

This paper introduces general ideas and some basic methods of the Bayesian probability theory applied to physics measurements. Our aim is to make the reader familiar, through examples rather than rigorous formalism, with concepts such as the following: model comparison (including the automatic Ockham's Razor filter provided by the Bayesian approach); parametric inference; quantification of the uncertainty about the value of physical quantities, also taking into account systematic effects; role of marginalization; posterior characterization; predictive distributions; hierarchical modelling and hyperparameters; Gaussian approximation of the posterior and recovery of conventional methods, especially maximum likelihood and chi-square fits under well-defined conditions; conjugate priors, transformation invariance and maximum entropy motivated priors; and Monte Carlo (MC) estimates of expectation, including a short introduction to Markov Chain MC methods.

## 1. Introduction

The last decades of the twentieth century have seen a considerable increase in the use of Bayesian methods in all fields of human activity that generally deal with uncertainty, including engineering, computer science, economics, medicine and even forensics (Kadane and Schum 1996). Bayesian networks (Pearl 1988, Cowell *et al* 1999) are used to diagrammatically represent uncertainty in expert systems or to construct artificial intelligence systems. Even venerable metrological associations, such as the International Organization for Standardization (ISO 1993), the Deutsches Institut für Normung (DIN 1996, 1999), and the USA National Institute of Standards and Technology (Taylor and Kuyatt 1994), have come to realize that Bayesian ideas are essential to provide general methods for quantifying uncertainty in measurements. A short account of the Bayesian upsurge can be found in Malakoff (1999). A search on the web for the keywords 'Bayesian', 'Bayesian network', or 'belief network' gives a dramatic impression of this 'revolution', not only in terms of improved methods, but more importantly in terms of reasoning. An overview of recent developments in Bayesian statistics may be found in the proceedings of the Valencia Conference series. The last published volume was Bernardo *et al* (1999), and the most recent conference was held in June 2002. Another series of workshops, under the title Maximum Entropy and Bayesian Methods, has focused more on applications in the physical sciences.

It is surprising that many physicists have been slow to adopt these 'new' ideas. There have been notable exceptions, of course, many of whom have contributed to the above-mentioned Maximum Entropy workshops. One reason to be surprised is that numerous great physicists and mathematicians have played important roles in developing probability theory. These 'new' ideas actually originated long ago with Bernoulli, Laplace, and Gauss, to mention just a few who contributed significantly to the development of physics, as well as to Bayesian thinking. So, while modern statisticians and mathematicians are developing powerful methods to apply to Bayesian analysis, most physicists, in their use and reasoning in statistics, still rely on twentieth century 'frequentist prescriptions' (D'Agostini 1999a, 2000).

We hope that this paper will help fill this gap by reviewing the advantages of using the Bayesian approach to address physics problems. We will emphasize the intuitive and practical aspects more than the theoretical ones. We will not try to cover all possible applications of Bayesian analysis in physics, but concentrate mainly on some basic applications that illustrate clearly the power of the method and how naturally it meshes with physicists' approach to their science.

The vocabulary, expressions, and examples have been chosen with the intent to correspond, as closely as possible, to the education that physicists receive in statistics, instead of a more rigorous approach that formal Bayesian statisticians might prefer. For example, we avoid many important theoretical concepts, like exchangeability, and do not attempt to prove the basic rules of probability. When we talk about 'random variables', we will in fact mean 'uncertain variables', instead of referring to the frequentist concept of 'randomness' *à la* von Mises (1957). This distinction will be clarified later.

In the past, presentations on Bayesian probability theory often started with criticisms of 'conventional', that is, frequentist, ideas, methods, and results. We shall keep criticisms and detailed comparisons of the results of different methods to a minimum. Readers interested in a critical review of conventional frequentist statistics will find a large literature, because most introductory books or reports on Bayesian analysis contain enough material on this matter (see Loredó 1990, Gelman *et al* 1995, Sivia 1997, Jaynes 1998, D'Agostini 2003 and references therein). Eloquent 'defenses of the Bayesian choice' can be found in Howson and Urbach (1993) and Robert (2001).

Some readers may wish to have references to unbiased comparisons of frequentist ideas and methods with Bayesian ones. To our knowledge, no such reports exist. Those who claim to be impartial are often frequentists who take some Bayesian results as if they were frequentist ‘prescriptions’, not caring whether all underlying hypotheses apply. For two prominent papers of this kind, see the articles by Efron (1986a) (with follow up discussions by Lindley (1986), Zellner (1986), and Efron (1986b)) and Cousins (1995). A recent, pragmatic comparison of frequentist and Bayesian confidence limits can be found in Zech (2002).

Despite its lack of widespread use in physics, and its complete absence in physics courses (D’Agostini 1999a), Bayesian data analysis is increasingly being employed in many areas of physics, for example, in astronomy (Babu and Feigelson 1992, 1997, Gregory and Loredo 1992, 1996, Bontekoe *et al* 1994, Gregory 1999), in geophysics (Glimm and Sharp 1999), in high-energy physics (D’Agostini and Degrossi 1999, Ciuchini *et al* 2001), in image reconstruction (Hanson 1993), in microscopy (Higdon and Yamamoto 2001), in quantum Monte Carlo (MC) (Gubernatis *et al* 1991), and in spectroscopy (Skilling 1992, Fischer *et al* 1997, 1998, 2000), just to mention a few articles written in the last decade. Other examples will be cited throughout this paper.

## 2. Uncertainty and probability

In the practice of science, we constantly find ourselves in a state of uncertainty. Uncertainty about the data that an experiment will yield. Uncertainty about the true value of a physical quantity, even after an experiment has been done. Uncertainty about model parameters, calibration constants, and other quantities that might influence the outcome of the experiment, and hence influence our conclusions about the quantities of interest, or the models that might have produced the observed results.

In general, we know through experience that not all the events that could happen, or all conceivable hypotheses, are equally likely. Let us consider the outcome of you measuring the temperature at the location where you are presently reading this paper, assuming you use a digital thermometer with one degree resolution (or you round the reading to the degree if you have a more precise instrument). There are some values that you expect the thermometer to display, others you expect less, and extremes you do not believe at all (some of them are simply excluded by the thermometer you are going to use). Given two events  $E_1$  and  $E_2$ , for example,  $E_1$ : ‘ $T = 22^\circ\text{C}$ ’ and  $E_2$ : ‘ $T = 33^\circ\text{C}$ ’, you might consider  $E_2$  much more probable than  $E_1$ , just meaning that you believe  $E_2$  more likely to be the case than  $E_1$ . We could use different expressions to mean exactly the same thing: you consider  $E_2$  more likely; you have more confidence in  $E_2$ ; having to choose between  $E_1$  and  $E_2$  to win a prize, you would promptly choose  $E_2$ ; having to classify with a number, which we shall denote by  $P$ , your degree of confidence on the two outcomes, you would write  $P(E_2) > P(E_1)$ ; and many others.

On the other hand, we might prefer to state the opposite, i.e.  $P(E_1) > P(E_2)$ , with the same meaning of symbols and referring exactly to the same events: what you are going to read at your location with your thermometer. The reason is simply that we do not share the same information. We do not know who you are and where you are at this very moment. You and we are uncertain about the same event, but in a different way. Values that might appear very probable to you now, appear quite improbable, though not impossible, to us.

In this example, we have introduced two crucial aspects of the Bayesian approach:

- (i) As it is used in everyday language, the term probability has the intuitive meaning of ‘the degree of belief that an event will occur’.

- (ii) Probability depends on our state of knowledge, which is usually different for different people. In other words, probability is unavoidably subjective.

At this point, you might find all of this quite natural, and wonder why these intuitive concepts go by the esoteric name 'Bayesian'. We agree! The fact is that the main thrust of statistics theory and practice during the twentieth century has been based on a different concept of probability, in which it is defined as the limit of the long-term relative frequency of the outcome of these events. It revolves around the theoretical notion of infinite ensembles of 'identical experiments'. Without entering into an unavoidably long critical discussion of the frequentist approach, we simply want to point out that in such a framework, there is no way to introduce the probability of hypotheses. All practical methods to overcome this deficiency yield misleading, and even absurd, conclusions (see D'Agostini (2003) for several examples and also for a justification of why frequentistic tests 'often work').

Instead, if we recover the intuitive concept of probability, we are able to talk in a natural way about the probability of any kind of event, or, extending the concept, of any proposition. In particular, the probability evaluation based on the relative frequency of similar events that occurred in the past is easily recovered in Bayesian theory, under precise conditions of validity (see section 5.3). Moreover, a simple theorem from probability theory, Bayes' theorem, which we shall see in the next section, allows us to update probabilities on the basis of new information. This inferential use of Bayes' theorem is only possible if probability is understood in terms of degree of belief. Therefore, the terms 'Bayesian' and 'based on subjective probability' are practically synonyms, and usually mean 'in contrast to the frequentist, or conventional, statistics'. The terms 'Bayesian' and 'subjective' should be considered transitional. In fact, there is already a tendency among many Bayesians to simply refer to 'probabilistic methods', and so on (Jeffreys 1961, de Finetti 1974, Jaynes 1998, Cowell *et al* 1999).

As mentioned above, Bayes' theorem plays a fundamental role in probability theory. This means that subjective probabilities of logically connected events are related to each other by mathematical rules. This important result can be summed up by saying, in practical terms, that 'degrees of belief follow the same grammar as abstract axiomatic probabilities'. Hence, all formal properties and theorems from probability theory follow.

Within the Bayesian school, there is no single way to derive the basic rules of probability (note that they are not simply taken as axioms in this approach). de Finetti's principle of coherence (de Finetti 1974) is considered the best guidance by many leading Bayesians (Bernardo and Smith 1994, O'Hagan 1994, Lad 1996, Coletti and Scozzafava 2002). See D'Agostini (2003) for an informal introduction to the concept of coherence, which in simple words can be outlined as follows. A person who evaluates probability values should be ready to accept bets in either direction, with odd ratios calculated from those values of probability. For example, an analyst who claims to be 50% confident on  $E$  should be aware that somebody could ask him to make a 1:1 bet on  $E$  or on  $\bar{E}$ . If he/she feels uneasy, it means that he/she does not consider the two events equally likely and the 50% was 'incoherent'.

Others, in particular practitioners close to the Jaynes' Maximum Entropy school (Jaynes 1957a,b), feel more at ease with Cox's logical consistency reasoning, requiring some consistency properties ('desiderata') between values of probability related to logically connected propositions (Cox 1946) (see also Jaynes (1998), Sivia (1997), Fröhner (2000), and especially Tribus (1969), for accurate derivations and a clear account of the meaning and role of information entropy in data analysis). An approach similar to Cox's is followed by Jeffreys (1961), another leading figure who has contributed a new vitality to the methods based on this 'new' point of view on probability. Note that Cox and Jeffreys were physicists. Remarkably, Schrödinger (1947a,b) also arrived at similar conclusions, though his definition

of event is closer to de Finetti's. Some short quotations from (Schrödinger 1947a) are in order. Definition of probability:

... a quantitative measure of the strength of our conjecture or anticipation, founded on the said knowledge, that the event comes true'. Subjective nature of probability: 'Since the knowledge may be different with different persons or with the same person at different times, they may anticipate the same event with more or less confidence, and thus different numerical probabilities may be attached to the same event'. Conditional probability: 'Thus whenever we speak loosely of 'the probability of an event', it is always to be understood: probability with regard to a certain given state of knowledge.

### 3. Rules of probability

We begin by stating some basic rules and general properties that form the 'grammar' of the probabilistic language that is used in Bayesian analysis. In this section, we review the rules of probability, starting with the rules for simple propositions. We will not provide rigorous derivations and will not address the foundational or philosophical aspects of probability theory. Moreover, following an 'eclectic' approach, which is common among Bayesian practitioners, we talk indifferently about probability of events, probability of hypotheses, or probability of propositions. Indeed, the last expression will often be favoured, understanding that it includes the others.

#### 3.1. Probability of simple propositions

Let us start by recalling the basic rules of probability for propositions or hypotheses. Let  $A$  and  $B$  be propositions, which can take on only two values, for example, true or false. The notation  $P(A)$  stands for the probability that  $A$  is true. The elementary rules of probability for simple propositions are

$$0 \leq P(A) \leq 1, \quad (1)$$

$$P(\Omega) = 1, \quad (2)$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B), \quad (3)$$

$$P(A \cap B) = P(A | B) P(B) = P(B | A) P(A), \quad (4)$$

where  $\Omega$  means tautology (a proposition that is certainly true). The construct  $A \cap B$  is true only when both  $A$  and  $B$  are true (logical AND), while  $A \cup B$  is true when at least one of the two propositions is true (logical OR).  $A \cap B$  is also written simply as ' $A, B$ ' or  $AB$ , and is also called a logical product, while  $A \cup B$  is also called a logical sum.  $P(A, B)$  is called the joint probability of  $A$  and  $B$ .  $P(A | B)$  is the probability of  $A$  under the condition that  $B$  is true. We often read it simply as 'the probability of  $A$ , given  $B$ '.

Equation (4) shows that the joint probability of two events can be decomposed into conditional probabilities in two different ways. Either of these ways is called the product rule. If the status of  $B$  does not change the probability of  $A$ , and the other way around, then  $A$  and  $B$  are said to be independent, probabilistically independent to be precise. In that case,  $P(A | B) = P(A)$ , and  $P(B | A) = P(B)$ , which, when inserted in equation (4), yields

$$P(A \cap B) = P(A) P(B) \iff \text{probabilistic independence.} \quad (5)$$

Equations (1)–(4) logically lead to other rules, which form the body of probability theory. For example, indicating the negation (or opposite) of  $A$  with  $\bar{A}$ , clearly  $A \cup \bar{A}$  is a tautology

$(A \cup \bar{A} = \Omega)$ , and  $A \cap \bar{A}$  is a contradiction ( $A \cap \bar{A} = \emptyset$ ). The symbol  $\emptyset$  stands for contradiction (a proposition that is certainly false). Hence, we obtain from equations (2) and (3)

$$P(A) + P(\bar{A}) = 1, \quad (6)$$

which says that proposition  $A$  is either true or not true.

### 3.2. Probability of complete classes

These formulae become more interesting when we consider a set of propositions  $H_j$  that all together form a tautology (i.e. they are exhaustive) and are mutually exclusive. Formally,

$$\bigcup_i H_j = \Omega, \quad (7)$$

$$H_j \cap H_k = \emptyset \quad \text{if } j \neq k. \quad (8)$$

When these conditions apply, the set  $\{H_j\}$  is said to form a complete class. The symbol  $H$  has been chosen because we shall soon interpret  $\{H_j\}$  as a set of hypotheses.

The first (trivial) property of a complete class is normalization, that is,

$$\sum_j P(H_j) = 1, \quad (9)$$

which is just an extension of equation (6) to a complete class containing more than just a single proposition and its negation.

For the complete class  $H$ , the generalizations of equation (6) and the use of equation (4) yield

$$P(A) = \sum_j P(A, H_j), \quad (10)$$

$$P(A) = \sum_j P(A | H_j) P(H_j). \quad (11)$$

Equation (10) is the basis of what is called marginalization, which will become particularly important when dealing with uncertain variables: the probability of  $A$  is obtained by the summation over all possible constituents contained in  $A$ . Hereafter, we avoid explicitly writing the limits of the summations, meaning that they extend over all elements of the class. The constituents are ' $A, H_j$ ', which, based on the complete class of hypotheses  $\{H\}$ , themselves form a complete class, which can easily be proved. Equation (11) shows that the probability of any proposition is given by a weighted average of all conditional probabilities, subject to hypotheses  $H_j$  forming a complete class, with the weight being the probability of the hypothesis.

In general, there are many ways to choose complete classes (like 'bases' in geometrical spaces). Let us denote the elements of a second complete class by  $E_i$ . The constituents are then formed by the elements  $(E_i, H_j)$  of the Cartesian product  $\{E\} \times \{H\}$ . Equations (10) and (11) then become the more general statements

$$P(E_i) = \sum_j P(E_i, H_j), \quad (12)$$

$$P(E_i) = \sum_j P(E_i | H_j) P(H_j), \quad (13)$$

and, symmetrically,

$$P(H_j) = \sum_i P(E_i, H_j), \quad (14)$$

$$P(H_j) = \sum_i P(H_j | E_i) P(E_i). \quad (15)$$

The reason we write these formulae both ways is to stress the symmetry of Bayesian reasoning with respect to classes  $\{E\}$  and  $\{H\}$ , though we shall soon associate them with observations (or events) and hypotheses, respectively.

### 3.3. Probability rules for uncertain variables

In analysing the data from physics experiments, we need to deal with measurements that are discrete or continuous in nature. Our aim is to make inferences about the models that we believe appropriately describe the physical situation, and/or, within a given model, to determine the values of relevant physics quantities. Thus, we need the probability rules that apply to uncertain variables, whether they are discrete or continuous. The rules for complete classes described in the preceding section clearly apply directly to discrete variables. With only slight changes, the same rules also apply to continuous variables because they may be thought of as a limit of discrete variables, as the interval between possible discrete values goes to zero.

For a discrete variable  $x$ , the expression  $p(x)$ , which is called a probability function, has the interpretation in terms of the probability of the proposition  $P(A)$ , where  $A$  is true when the value of the variable is equal to  $x$ . In the case of continuous variables, we use the same notation, but with the meaning of a probability density function (pdf). So  $p(x) dx$ , in terms of a proposition, is the probability  $P(A)$ , where  $A$  is true when the value of the variable lies in the range  $x$  to  $x + dx$ . In general, the meaning is clear from the context; otherwise it should be stated. Probabilities involving more than one variable, like  $p(x, y)$ , have the meaning of the probability of a logical product; they are usually called joint probabilities.

Table 1 summarizes useful formulae for discrete and continuous variables. The interpretation and use of these relations in Bayesian inference will be illustrated in the following sections.

**Table 1.** Some definitions and properties of probability functions for values of a discrete variable  $x_i$  and probability density functions for continuous variables  $x$ . All summations and integrals are understood to extend over the full range of possibilities of the variable. Note that the expectation of the variable is also called the expected value (sometimes expectation value), average, and mean. The square root of the variance is the standard deviation  $\sigma$ .

	Discrete variables	Continuous variables
Probability	$P(X = x_i) = p(x_i)$	$dP_{[x \leq X \leq x+dx]} = p(x) dx$
Normalization <sup>a</sup>	$\sum_i p(x_i) = 1$	$\int p(x) dx = 1$
Expectation of $f(X)$	$E[f(X)] = \sum_i f(x_i) p(x_i)$	$E[f(X)] = \int f(x) p(x) dx$
Expected value	$E(X) = \sum_i x_i p(x_i)$	$E(X) = \int x p(x) dx$
Moment of order $r$	$M_r(X) = \sum_i x_i^r p(x_i)$	$M_r(X) = \int x^r p(x) dx$
Variance	$\sigma^2 = \sum_i [x_i - E(X)]^2 p(x_i)$	$\sigma^2 = \int [x - E(X)]^2 p(x) dx$
Product rule	$p(x_i, y_j) = p(x_i   y_j) p(y_j)$	$p(x, y) = p(x   y) p(y)$
Independence	$p(x_i, y_j) = p(x_i) p(y_j)$	$p(x, y) = p(x) p(y)$
Marginalization	$\sum_j p(x_i, y_j) = p(x_i)$	$\int p(x, y) dy = p(x)$
Decomposition	$p(x_i) = \sum_j p(x_i   y_j) p(y_j)$	$p(x) = \int p(x   y) p(y) dy$
Bayes' theorem	$p(x_j   y_i) = \frac{p(y_i   x_j) p(x_j)}{\sum_j p(y_i   x_j) p(x_j)}$	$p(x   y) = \frac{p(y   x) p(x)}{\int p(y   x) p(x) dx}$
Likelihood	$\mathcal{L}(x_j; y_i) = p(y_i   x_j)$	$\mathcal{L}(x; y) = p(y   x)$

<sup>a</sup> A function  $p(x)$  such that  $\sum_i p(x_i) = \infty$ , or  $\int p(x) dx = \infty$ , is called improper. Improper functions are often used to describe relative beliefs about the possible values of a variable.

#### 4. Bayesian inference for simple problems

We introduce the basic concepts of Bayesian inference by considering some simple problems. The aim is to illustrate some of the notions that form the foundation of Bayesian reasoning.

##### 4.1. Background information

As we think about drawing conclusions about the physical world, we come to realize that everything we do is based on what we know about the world. Conclusions about hypotheses will be based on our general background knowledge. To emphasize the dependence of probability on the state of background information, which we designate by  $I$ , we will make it explicit by writing  $P(E | I)$ , rather than simply  $P(E)$ . (Note that, in general,  $P(A | I_1) \neq P(A | I_2)$ , if  $I_1$  and  $I_2$  are different states of information.) For example, equation (4) should be more precisely written as

$$P(A \cap B | I) = P(A | B \cap I) P(B | I) = P(B | A \cap I) P(A | I), \quad (16)$$

or alternatively as

$$P(A, B | I) = P(A | B, I) P(B | I) = P(B | A, I) P(A | I). \quad (17)$$

We have explicitly included  $I$  as part of the conditional as a reminder that any probability relation is valid only under the same state of background information.

##### 4.2. Bayes' theorem

Formally, Bayes' theorem follows from the symmetry of  $P(A, B)$  expressed by equation (17). In terms of  $E_i$  and  $H_j$  belonging to two different complete classes, equation (17) yields

$$\frac{P(H_j | E_i, I)}{P(H_j | I)} = \frac{P(E_i | H_j, I)}{P(E_i | I)}. \quad (18)$$

This equation says that the new condition  $E_i$  alters our belief in  $H_j$  by the same updating factor by which the condition  $H_j$  alters our belief about  $E_i$ . Rearrangement yields Bayes' theorem:

$$P(H_j | E_i, I) = \frac{P(E_i | H_j, I) P(H_j | I)}{P(E_i | I)}. \quad (19)$$

We have obtained a logical rule to update our beliefs on the basis of new conditions. Note that, though Bayes' theorem is a direct consequence of the basic rules of axiomatic probability theory, its updating power can only be fully exploited if we can treat on the same basis expressions concerning hypotheses and observations, causes and effects, models and data.

In most practical cases, the evaluation of  $P(E_i | I)$  can be quite difficult, while determining the conditional probability  $P(E_i | H_j, I)$  might be easier. For example, think of  $E_i$  as the probability of observing a particular event topology in a particle physics experiment, compared with the probability of the same thing given a value of the hypothesized particle mass ( $H_j$ ), a given detector, background conditions, etc. Therefore, it is convenient to rewrite  $P(E_i | I)$  in equation (19) in terms of the quantities in the numerator, using equation (13), to obtain

$$P(H_j | E_i, I) = \frac{P(E_i | H_j, I) P(H_j | I)}{\sum_j P(E_i | H_j, I) P(H_j | I)}, \quad (20)$$

which is the better-known form of Bayes' theorem. Written this way, it becomes evident that the denominator of the right-hand side of equation (20) is just a normalization factor and we can focus on just the numerator:

$$P(H_j | E_i, I) \propto P(E_i | H_j, I) P(H_j | I). \quad (21)$$



Saying the same thing in words,

$$\text{posterior} \propto \text{likelihood} \times \text{prior}, \quad (22)$$

where the posterior (or final state) stands for the probability of  $H_j$ , based on the new observation  $E_i$ , relative to the prior (or initial) probability. (Prior probabilities are often indicated by  $P_0$ .) The conditional probability  $P(E_i | H_j)$  is called the likelihood. It is literally the probability of the observation  $E_i$  given the specific hypothesis  $H_j$ . The term likelihood can lead to some confusion, because it is often misunderstood to mean ‘the likelihood that  $E_i$  comes from  $H_j$ ’. However, this name implies that we consider  $P(E_i | H_j)$  as a mathematical function of  $H_j$  for a fixed  $E_i$  and in that framework it is usually written as  $\mathcal{L}(H_j; E_i)$  to emphasize the functionality. We caution the reader that one sometimes even finds the notation  $\mathcal{L}(E_i | H_j)$  to indicate exactly  $P(E_i | H_j)$ .

#### 4.3. Inference for simple hypotheses

Making use of formulae (20) or (21), we can easily solve many classical problems involving inference when many hypotheses can produce the same single effect. Consider the case of interpreting the results of a test for the HIV virus applied to a randomly chosen European. Clinical tests are very seldom perfect. Suppose that the test accurately detects infection, but has a false-positive rate of 0.2%:

$$P(\text{Positive} | \text{Infected}) = 1 \quad \text{and} \quad P(\text{Positive} | \overline{\text{Infected}}) = 0.2\%.$$

If the test is positive, can we conclude that the particular person is infected with a probability of 99.8% because the test has only a 0.2% chance of mistake? Certainly not! This kind of mistake is often made by those who are not used to Bayesian reasoning, including scientists who make inferences in their own field of expertise. The correct answer depends on what else we know about the person tested, that is, the background information. Thus, we have to consider the incidence of the HIV virus in Europe, and possibly, information about the lifestyle of the individual (for details, see D’Agostini (2003)).

To better understand the updating mechanism, let us take the ratio of equation (20) for two hypotheses  $H_j$  and  $H_k$ :

$$\frac{P(H_j | E_i, I)}{P(H_k | E_i, I)} = \frac{P(E_i | H_j, I) P(H_j | I)}{P(E_i | H_k, I) P(H_k | I)}, \quad (23)$$

where the sums in the denominators of equation (20) cancel. It is convenient to interpret the ratio of probabilities, given the same condition, as betting odds. This is best done formally in the de Finetti approach, but the basic idea is what everyone is used to: the amount of money that one is willing to bet on an event is proportional to the degree to which one expects that event will happen. Equation (23) tells us that, when new information is available, the initial odds are updated by the ratio of the likelihoods  $P(E_i | H_j, I)/P(E_i | H_k, I)$ , which is known as the Bayes factor.

In the case of the HIV test, the initial odds for an arbitrarily chosen European being infected  $P(H_j | I)/P(H_k | I)$  are so small that we need a very high Bayes’ factor to be reasonably certain that, when the test is positive, the person is really infected. With the numbers used in this example, the Bayes factor is  $500 = 1/0.002$ . For example, if we take for the prior  $P_0(\text{Infected})/P_0(\overline{\text{Infected}}) = 1/1000$ , the Bayes’ factor changes these odds to  $500/1000 = \frac{1}{2}$ , or equivalently, the probability that the person is infected would be  $\frac{1}{3}$ , quite different from the 99.8% answer usually prompted by those who have a standard statistical education. This example can be translated straightforwardly to physical problems, like particle identification in the analysis of Cherenkov detector data, as done, for example, in D’Agostini (2003).

## 5. Inferring numerical values of physics quantities—general ideas and basic examples

In physics, we are concerned about models ('theories') and the numerical values of physical quantities related to them. Models and the values of quantities are, generally speaking, the hypothesis we want to infer, given the observations. In the previous section we learned how to deal with simple hypotheses—'simple' in the sense that they do not depend on internal parameters.

On the other hand, in many applications we have strong beliefs about what model to use to interpret the measurements. Thus, we focus our attention on the model parameters that we consider as uncertain variables that we want to infer. The method that deals with these applications is usually referred to as parametric inference, and it will be shown with examples in this section. In our models, the value of the relevant physical quantities are usually described in terms of a continuous uncertain variable. Bayes' theorem, properly extended to uncertain quantities (see table 1), plays a central role in this inference process.

A more complicated case is when we are also uncertain about the model (and each possible model has its own set of parameters, usually associated with different physics quantities). We shall analyse this problem in section 7.

### 5.1. Bayesian inference on uncertain variables and posterior characterization

We start here with a few one-dimensional problems involving simple models that often occur in data analysis. These examples will be used to illustrate some of the most important Bayesian concepts. Let us first introduce briefly the structure of the Bayes' theorem in the form convenient to our purpose, as a straightforward extension of what was seen in section 4.2.

$$p(\theta | d, I) = \frac{p(d | \theta, I) p(\theta | I)}{\int p(d | \theta, I) p(\theta | I) d\theta}, \quad (24)$$

$\theta$  is the generic name of the parameter (used hereafter, unless the models have traditional symbols for their parameters) and  $d$  is the data point.  $p(\theta | I)$  is the prior,  $p(\theta | d, I)$  the posterior, and  $p(d | \theta, I)$  the likelihood. Also, in this case, the likelihood is often written as  $\mathcal{L}(\theta; d) = p(d | \theta, I)$ , and the words of caution expressed in section 4.2 apply here too. Note, moreover, that, while  $p(d | \theta, I)$  is a properly normalized pdf,  $\mathcal{L}(\theta; d)$  does not have a pdf meaning in the variable  $\theta$ . Hence, the integral of  $\mathcal{L}(\theta; d)$  over  $\theta$  is only accidentally equal to unity. The denominator on the right-hand side of equation (24) is called the evidence and, while in the parametric inference discussed here it is just a trivial normalization factor, its value becomes important for model comparison (see section 7).

Posterior probability distributions provide the full description of our state of knowledge about the value of the quantity. In fact, they allow us to calculate all probability intervals of interest. Such intervals are also called credible intervals (at a specified level of probability, for example, 95%) or confidence intervals (at a specified level of 'confidence', i.e. of probability). However, the latter expression could be confused with frequentistic 'confidence intervals', which are not probabilistic statements about uncertain variables (D'Agostini 2003).

It is often desirable to characterize the distribution in terms of a few numbers, for example, mean value (arithmetic average) of the posterior, or its most probable value (the mode) of the posterior, also known as the maximum *a posteriori* (MAP) estimate. The spread of the distribution is often described in terms of its standard deviation (square root of the variance). It is useful to associate the terms mean value and standard deviation with the more inferential terms expected value, or simply expectation (value), indicated by  $E(\cdot)$ , and standard uncertainty (ISO 1993), indicated by  $\sigma(\cdot)$ , where the argument is the uncertain variable of interest. This will be our standard method of reporting the result of inference in a quantitative way, although

we emphasize that the full answer is given by the posterior distribution, and reporting only these summaries in the case of complex distributions (e.g. multimodal and/or asymmetrical pdfs) can be misleading, because people tend to think of a Gaussian model if no further information is provided.

5.2. Gaussian model

Let us start with a classical example, in which the response signal  $d$  from a detector is described by a Gaussian error function around the true value  $\mu$  with a standard deviation  $\sigma$ , which is assumed to be exactly known. This model is the best known among physicists and, indeed, the Gaussian pdf is also known as normal because it is often assumed that errors are ‘normally’ distributed according to this function. Applying Bayes’ theorem for continuous variables (see table 1), from the likelihood

$$p(d | \mu, I) = \frac{1}{\sqrt{2\pi} \sigma} \exp \left[ -\frac{(d - \mu)^2}{2 \sigma^2} \right], \tag{25}$$

we get for  $\mu$

$$p(\mu | d, I) = \frac{(1/\sqrt{2\pi} \sigma) \exp[-(d - \mu)^2/2 \sigma^2] p(\mu | I)}{\int_{-\infty}^{+\infty} (1/\sqrt{2\pi} \sigma) \exp[-(d - \mu)^2/2 \sigma^2] p(\mu | I) d\mu}. \tag{26}$$

Considering all values of  $\mu$  equally likely over a very large interval, we can model the prior  $p(\mu | I)$  with a constant, which simplifies in equation (26), yielding

$$p(\mu | d, I) = \frac{1}{\sqrt{2\pi} \sigma} \exp \left[ -\frac{(\mu - d)^2}{2 \sigma^2} \right]. \tag{27}$$

Expectation and standard deviation of the posterior distribution are  $E(\mu) = d$  and  $\sigma(\mu) = \sigma$ , respectively. This particular result corresponds to what is often done intuitively in practice. But one has to pay attention to the assumed conditions under which the result is logically valid: Gaussian likelihood and uniform prior. Moreover, we can speak about the probability of true values only in the subjective sense. It is recognized that physicists, and scientists in general, are highly confused about this point (D’Agostini 1999a).

A noteworthy case of a prior for which the naive inversion gives paradoxical results is when the value of a quantity is constrained to be in the ‘physical region’, for example,  $\mu \geq 0$ , while  $d$  falls outside it (or is at its edge). The simplest prior that cures the problem is a step function  $\theta(\mu)$ , and the result is equivalent to simply renormalizing the pdf in the physical region (this result corresponds to a ‘prescription’ sometimes used by practitioners with a frequentist background when they encounter this kind of problem).

Another interesting case is when the prior knowledge can be modelled with a Gaussian function, for example, describing our knowledge from a previous inference:

$$p(\mu | \mu_0, \sigma_0, I) = \frac{1}{\sqrt{2\pi} \sigma_0} \exp \left[ -\frac{(\mu - \mu_0)^2}{2 \sigma_0^2} \right]. \tag{28}$$

Inserting equation (28) into equation (26), we get

$$p(\mu | d, \mu_0, \sigma_0, I) = \frac{1}{\sqrt{2\pi} \sigma_1} \exp \left[ -\frac{(\mu - \mu_1)^2}{2 \sigma_1^2} \right], \tag{29}$$

where

$$\mu_1 = E(\mu) = \frac{d/\sigma^2 + \mu_0/\sigma_0^2}{1/\sigma^2 + 1/\sigma_0^2}, \tag{30}$$

$$\mu_1 = E(\mu) = \frac{\sigma_0^2}{\sigma^2 + \sigma_0^2} d + \frac{\sigma^2}{\sigma^2 + \sigma_0^2} \mu_0 = \frac{\sigma_1^2}{\sigma^2} d + \frac{\sigma_1^2}{\sigma_0^2} \mu_0, \quad (31)$$

$$\sigma_1^2 = \text{Var}(\mu) = (\sigma_0^{-2} + \sigma^{-2})^{-1}. \quad (32)$$

We can then see that the case  $p(\mu | I) = \text{constant}$  corresponds to the limit of a Gaussian prior with very large  $\sigma_0$  and finite  $\mu_0$ . The formula for the expected value combining previous knowledge and present experimental information has been written in several ways in equation (31).

Another enlightening way of writing equation (30) is to consider  $\mu_0$  and  $\mu_1$  the estimates of  $\mu$  at times  $t_0$  and  $t_1$ , respectively, before and after the observation  $d$  happened at time  $t_1$ . Indicating the estimates at different times by  $\hat{\mu}(t)$ , we can rewrite equation (30) as

$$\begin{aligned} \hat{\mu}(t_1) &= \frac{\sigma_\mu^2(t_0)}{\sigma_d^2(t_1) + \sigma_\mu^2(t_0)} d(t_1) + \frac{\sigma_d^2(t_1)}{\sigma_d^2(t_1) + \sigma_\mu^2(t_0)} \hat{\mu}(t_0) \\ &= \hat{\mu}(t_0) + \frac{\sigma_\mu^2(t_0)}{\sigma_d^2(t_1) + \sigma_\mu^2(t_0)} [d(t_1) - \hat{\mu}(t_0)], \end{aligned} \quad (33)$$

$$\hat{\mu}(t_1) = \hat{\mu}(t_0) + K(t_1) [d(t_1) - \hat{\mu}(t_0)], \quad (34)$$

$$\sigma_\mu^2(t_1) = \sigma_\mu^2(t_0) - K(t_1) \sigma_\mu^2(t_0), \quad (35)$$

where

$$K(t_1) = \frac{\sigma_\mu^2(t_0)}{\sigma_d^2(t_1) + \sigma_\mu^2(t_0)}. \quad (36)$$

Indeed, we have given equation (30) the structure of a Kalman filter (Kalman 1960). The new observation 'corrects' the estimate by a quantity given by the innovation (or residual)  $[d(t_1) - \hat{\mu}(t_0)]$  times the blending factor (or gain)  $K(t_1)$ . For an introduction to the Kalman filter and its probabilistic origin (see Maybeck (1979), Welch and Bishop (2002)).

As equations (31)–(35) show, new experimental information reduces the uncertainty. But this is true as long the previous information and the observation are somewhat consistent. If we are, for several reasons, sceptical about the model which yields the combination rule (31) and (32), we need to remodel the problem and introduce possible systematic errors or underestimations of the quoted standard deviations, as done, for example, in Press (1997), Dose and von der Linden (1999), D'Agostini (1999b), Fröhner (2000).

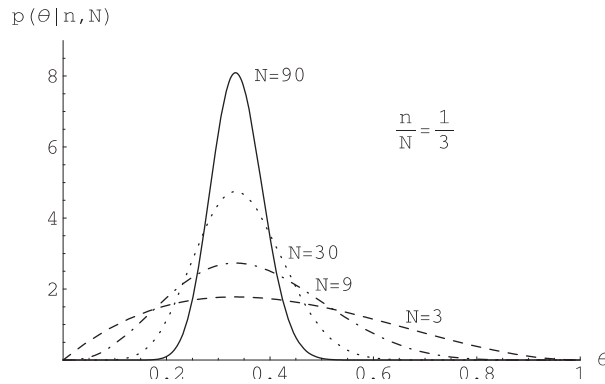
### 5.3. Binomial model

In a large class of experiments, the observations consist of counts, that is, a number of things (events, occurrences, etc). In many processes of physics interests the resulting number of counts is described probabilistically by a binomial or a Poisson model. For example, we would like to draw an inference about the efficiency of a detector, a branching ratio in a particle decay, or a rate from a measured number of counts in a given interval of time.

The binomial distribution describes the probability of randomly obtaining  $n$  events ('successes') in  $N$  independent trials, in each of which we assume the same probability  $\theta$  that the event will happen. The probability function is

$$p(n | \theta, N) = \binom{N}{n} \theta^n (1 - \theta)^{N-n}, \quad (37)$$

where the leading factor is the well-known binomial coefficient, namely  $N!/n!(N-n)!$ . We wish to infer  $\theta$  from an observed number of counts  $n$  in  $N$  trials. Incidentally, that was



**Figure 1.** Posterior probability density function of the binomial parameter  $\theta$ , having observed  $n$  successes in  $N$  trials.

the ‘problem in the doctrine of chances’ originally treated by Bayes (1763), reproduced, for example, in Press (1997). Assuming a uniform prior for  $\theta$ , by Bayes’ theorem the posterior distribution for  $\theta$  is proportional to the likelihood, given by equation (37):

$$p(\theta | n, N, I) = \frac{\theta^n (1 - \theta)^{N-n}}{\int_0^1 \theta^n (1 - \theta)^{N-n} d\theta}, \tag{38}$$

$$p(\theta | n, N, I) = \frac{(N + 1)!}{n! (N - n)!} \theta^n (1 - \theta)^{N-n}. \tag{39}$$

Some examples of this distribution for various values of  $n$  and  $N$  are shown in figure 1. The expectation, variance, and mode of this distribution are

$$E(\theta) = \frac{n + 1}{N + 2}, \tag{40}$$

$$\sigma^2(\theta) = \frac{(n + 1)(N - n + 1)}{(N + 3)(N + 2)^2} = \frac{E(\theta) (1 - E(\theta))}{N + 3}, \tag{41}$$

$$\theta_m = \frac{n}{N}, \tag{42}$$

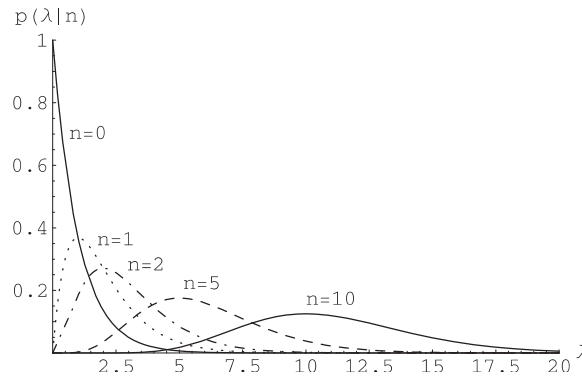
where the mode has been indicated by  $\theta_m$ . Equation (40) is known as the Laplace formula. For large values of  $N$  and  $0 \ll n \ll N$  the expectation of  $\theta$  tends to  $\theta_m$ , and  $p(\theta)$  becomes approximately Gaussian. This result is nothing but a reflection of the well-known asymptotic Gaussian behaviour of  $p(n | \theta, N)$ . For large  $N$  the uncertainty about  $\theta$  goes like  $1/\sqrt{N}$ . Asymptotically, we are practically certain that  $\theta$  is equal to the relative frequency of that class of events observed in the past. This is how the frequency based evaluation of probability is promptly recovered in the Bayesian approach, under well-defined assumptions.

5.4. Poisson model

The Poisson distribution gives the probability of observing  $n$  counts in a fixed time interval, when the expectation of the number of counts to be observed is  $\lambda$ :

$$p(n | \lambda) = \frac{\lambda^n e^{-\lambda}}{n!}. \tag{43}$$

The inverse problem is to infer  $\lambda$  from  $n$  counts observed. (Note that what matters physically is the rate  $r = \lambda/\Delta T$ , where  $\Delta T$  is the observation time.) Applying Bayes’ theorem and using



**Figure 2.** The posterior distribution for the Poisson parameter  $\lambda$ , when  $n$  counts are observed in an experiment.

a uniform prior  $p(\lambda | I)$  for  $\lambda$ , we get

$$p(\lambda | n, I) = \frac{\lambda^n e^{-\lambda}/n!}{\int_0^{\infty} (\lambda^n e^{-\lambda}/n!) d\lambda} = \frac{\lambda^n e^{-\lambda}}{n!}. \quad (44)$$

Some examples of this distribution for various values of  $n$  are shown in figure 2. As for the Gaussian model, the same mathematical expression holds for the likelihood, but with interchanged role of variable and parameter. The expectation and variance of  $\lambda$  are both equal to  $n + 1$ , while the most probable value is  $\lambda_m = n$ . For large  $n$ , the extra '+1' (due to the asymmetry of the prior with respect to  $\lambda = 0$ ) can be ignored and we have  $E(\lambda) = \sigma^2(\lambda) \approx n$  and, once again, the uncertainty about  $\lambda$  follows a Gaussian model. The relative uncertainty in  $\lambda$  decreases as  $1/\sqrt{n}$ .

When the observed value of  $n$  is zero, equation (44) yields  $p(\lambda | n = 0) = e^{-\lambda}$ , giving a maximum of belief at zero, but an exponential tail towards large values of  $\lambda$ . The expected value and standard deviation of  $\lambda$  are both equal to 1. The 95% probabilistic upper bound of  $\lambda$  is at  $\lambda_{95\%UB} = 3$ , as it can be easily calculated solving the equation  $\int_0^{\lambda_{95\%UB}} p(\lambda | n = 0) d\lambda = 0.95$ . Note, also, that this result depends on the choice of prior, though Astone and D'Agostini (1999) have shown that the upper bound is insensitive to the exact form of the prior, if the prior models somehow what they call 'positive attitude of rational scientists' (the prior does not have to be in contradiction with what one could actually observe, given the detector sensitivity). In particular, they show that a uniform prior is a good practical choice to model this attitude. On the other hand, talking about 'objective' probabilistic upper/lower limits makes no sense, as discussed in detail and with examples in the cited paper: one can at most speak about conventionally defined non-probabilistic sensitivity bounds, which separate the measurement region from that in which experimental sensitivity is lost (Astone and D'Agostini 1999, D'Agostini 2000, Astone *et al* 2002).

### 5.5. Inference from a data set and sequential use of the Bayes formula

In the elementary examples shown above, the inference has been made from a single data point  $d$ . If we have a set of observations (data), indicated by  $\mathbf{d}$ , we just need to insert in the Bayes formula the likelihood  $p(\mathbf{d} | \theta, I)$ , where this expression indicates a multi-dimensional joint pdf.

Note that we could think of inferring  $\theta$  on the basis of each newly observed datum  $d_i$ . After one observation

$$p(\theta | d_1, I) \propto p(d_1 | \theta, I) p(\theta | I), \quad (45)$$

and after the second

$$p(\theta | d_1, d_2, I) \propto p(d_2 | \theta, d_1, I) p(\theta | d_1, I), \quad (46)$$

$$p(\theta | d_1, d_2, I) \propto p(d_2 | \theta, d_1, I) p(d_1 | \theta, I) p(\theta | I). \quad (47)$$

We have written equation (47) in such a way that the dependence between observables can be accommodated. From the product rule in table 1, we can rewrite equation (47) as

$$p(\theta | d_1, d_2, I) \propto p(d_1, d_2 | \theta, I) p(\theta | I). \quad (48)$$

Comparing this equation with (47) we see that the sequential inference gives exactly the same result of a single inference that properly takes into account all available information. This is an important result of the Bayesian approach.

The extension to many variables is straightforward and we obtain

$$p(\boldsymbol{\theta} | \mathbf{d}, I) \propto p(\mathbf{d} | \boldsymbol{\theta}, I) p(\boldsymbol{\theta} | I). \quad (49)$$

Furthermore, when the  $d_i$  are independent, we get for the likelihood

$$p(\mathbf{d} | \boldsymbol{\theta}, I) = \prod_i p(d_i | \boldsymbol{\theta}, I), \quad (50)$$

$$\mathcal{L}(\boldsymbol{\theta}; \mathbf{d}) = \prod_i \mathcal{L}(\boldsymbol{\theta}; d_i), \quad (51)$$

that is, the combined likelihood is given by the product of the individual likelihoods.

### 5.6. Multi-dimensional case—inferring $\mu$ and $\sigma$ of a Gaussian

So far we have only inferred one parameter of a model. The extension to many parameters is straightforward. Calling  $\boldsymbol{\theta}$  the set of parameters and  $\mathbf{d}$  the data, Bayes' theorem becomes

$$p(\boldsymbol{\theta} | \mathbf{d}, I) = \frac{p(\mathbf{d} | \boldsymbol{\theta}, I) p(\boldsymbol{\theta} | I)}{\int p(\mathbf{d} | \boldsymbol{\theta}, I) p(\boldsymbol{\theta} | I) d\boldsymbol{\theta}}. \quad (52)$$

Equation (52) gives the posterior for the full parameter vector  $\boldsymbol{\theta}$ . Marginalization (see table 1) allows calculation of the probability distribution for a single parameter, for example,  $p(\theta_i | \mathbf{d}, I)$ , by integrating over the remaining parameters. The marginal distribution  $p(\theta_i | \mathbf{d}, I)$  is then the complete result of the Bayesian inference on the parameter  $\theta_i$ . Though the characterization of the marginal is done in the usual way described in section 5.1, there is often an interest in summarizing some characters of the multi-dimensional posterior that are unavoidably lost in the marginalization (imagine marginalization as a kind of geometrical projection). Useful quantities are the covariance between parameters  $\theta_i$  and  $\theta_j$ , defined as

$$\text{Cov}(\theta_i, \theta_j) = E[(\theta_i - E[\theta_i])(\theta_j - E[\theta_j])]. \quad (53)$$

As is well known, quantities that give a more intuitive idea of what is going on are the correlation coefficients, defined as  $\rho(\theta_i, \theta_j) = \text{Cov}(\theta_i, \theta_j) / \sigma(\theta_i)\sigma(\theta_j)$ . Variances and covariances form the covariance matrix  $\mathbf{V}(\boldsymbol{\theta})$ , with  $V_{ii} = \text{Var}(\theta_i)$  and  $V_{ij} = \text{Cov}(\theta_i, \theta_j)$ . We recall also that convenient formulae to calculate variances and covariances are obtained from the expectation of the products  $\theta_i\theta_j$ , together with the expectations of the parameters:

$$V_{ij} = E(\theta_i\theta_j) - E(\theta_i) E(\theta_j). \quad (54)$$

As a first example of a multi-dimensional distribution from a data set, we can think, again, of the inference of the parameter  $\mu$  of a Gaussian distribution, but in the case when  $\sigma$  is also unknown and needs to be determined by the data. From equations (52), (50), and (25), with  $\theta_1 = \mu$  and  $\theta_2 = \sigma$  and neglecting overall normalization, we obtain

$$p(\mu, \sigma | \mathbf{d}, I) \propto \sigma^{-n} \exp\left[-\frac{\sum_{i=1}^n (d_i - \mu)^2}{2\sigma^2}\right] p(\mu, \sigma | I), \quad (55)$$

$$p(\mu | \mathbf{d}, I) = \int p(\mu, \sigma | \mathbf{d}, I) d\sigma, \quad (56)$$

$$p(\sigma | \mathbf{d}, I) = \int p(\mu, \sigma | \mathbf{d}, I) d\mu. \quad (57)$$

The closed form of equations (56) and (57) depends on the prior and, perhaps, for the most realistic choice of  $p(\mu, \sigma | I)$ , such a compact solution does not exist. But this is not an essential issue, given the present computational power. (For example, the shape of  $p(\mu, \sigma | I)$  can be easily inspected by a modern graphical tool.) We would like to stress here the conceptual simplicity of the Bayesian solution to the problem. (In the case in which the data set contains more than a dozen observations, a flat  $p(\mu, \sigma | I)$ , with the constraint  $\sigma > 0$ , can be considered a good practical choice.)

### 5.7. Predictive distributions

A related problem is to ‘infer’ what an experiment will observe given our best knowledge of the underlying theory and its parameters. Infer is within quotes because the term is usually used for model and parameters, rather than for observations. In this case people prefer to speak about prediction (or prevision). But we recall that in Bayesian reasoning there is conceptual symmetry between the uncertain quantities that enter the problem. Probability density functions describing as yet unobserved events are referred to as predictive distributions. There is a conceptual difference with the likelihood, which also gives a probability of observation, but under different hypotheses, as the following example clarifies.

Given  $\mu$  and  $\sigma$ , and assuming a Gaussian model, our uncertainty about a ‘future’  $d_f$  is described by the Gaussian pdf equation (25) with  $d = d_f$ . But this holds only under that particular hypothesis for  $\mu$  and  $\sigma$ , while, in general, we are also uncertain about these values. Applying the decomposition formula (table 1) we get:

$$p(d_f | I) = \int p(d_f | \mu, \sigma, I) p(\mu, \sigma | I) d\mu d\sigma. \quad (58)$$

Again, the integral might be technically difficult, but the solution is conceptually simple. Note that, though the decomposition formula is a general result of probability theory, it can be applied to this problem only in the subjective approach.

An analytically easy, insightful case is that of experiments with well-known  $\sigma$ . Given a past observation  $d_p$  and a vague prior,  $p(\mu | d_p, I)$  is Gaussian around  $d_p$  with variance  $\sigma_p^2$  (note that, with respect to  $p(\mu, \sigma | I)$  of equation (58), it has been made explicit that  $p(\mu)$  depends on  $d_p$ ).  $p(d_f | \mu)$  is Gaussian around  $\mu$  with variance  $\sigma_f^2$ . Finally, we get

$$p(d_f | d_p, I) = \int p(d_f | \mu, I) p(\mu | d_p, I) d\mu, \quad (59)$$

$$p(d_f | d_p, I) = \frac{1}{\sqrt{2\pi} \sqrt{\sigma_p^2 + \sigma_f^2}} \exp\left[-\frac{(d_f - d_p)^2}{2(\sigma_p^2 + \sigma_f^2)}\right]. \quad (60)$$



Note that in equation (59), independence of the observations, given  $\mu$ , was assumed, in the sense that  $p(d_p | \mu, I)$  and  $p(d_f | \mu, I)$  are stochastically independent: if  $\mu$  is precisely known, the past observation does not change our beliefs about what we can observe next. However, when  $\mu$  is not precisely known,  $d_p$  and  $d_f$  become dependent on each other, which is what equation (59) says.

### 5.8. Hierarchical modelling and hyperparameters

As we have seen in the previous section, it is often desirable to include in a probabilistic model one's uncertainty regarding various aspects of a pdf. This is a natural feature of Bayesian methods, due to the uniform approach used to deal with uncertainty and from which powerful analytical tools are derived. This kind of modelling is called hierarchical because the characteristics of one pdf are controlled by another pdf. All uncertain parameters on which the pdf depends are called hyperparameters. An example of the use of the hyperparameter is described in section 8.3 in which the prior to infer  $\theta$  in a binomial model is shown to be controlled by the parameters of a beta distribution.

As an example of practical importance, think of the combination of experimental results in the presence of outliers, i.e. of data points which are somehow in mutual disagreement. In this case the combination rule given by equations (30)–(32), extended to many data points, produces unacceptable conclusions. A way of solving this problem (Dose and von der Linden 1999, D'Agostini 1999b) is to model a scepticism about the quoted standard deviations of the experiments, introducing a pdf  $f(r)$ , where  $r$  is a rescaling factor of the standard deviation. In this way the  $\sigma$  that enter the right-hand side of equations (30)–(32) are hyperparameters of the problem. An alternative approach, also based on hierarchical modelling, is shown in Fröhner (2000) (for a more complete introduction to the subject see, e.g. Gelman *et al* (1995)).

### 5.9. From Bayesian inference to maximum-likelihood and minimum chi-squared model fitting

Let us continue with the case in which we know so little about appropriate values of the parameters that a uniform distribution is a practical choice for the prior. Equation (52) becomes

$$p(\theta | \mathbf{d}, I) \propto p(\mathbf{d} | \theta, I) p_0(\theta, I) \propto p(\mathbf{d} | \theta, I) = \mathcal{L}(\theta; \mathbf{d}), \quad (61)$$

where, we recall, the likelihood  $\mathcal{L}(\theta; \mathbf{d})$  is seen as a mathematical function of  $\theta$ , with parameters  $\mathbf{d}$ .

The set of  $\theta$  that is most likely is that which maximizes  $\mathcal{L}(\theta; \mathbf{d})$ , a result known as the maximum likelihood principle. Here, it has been obtained again as a special case of a more general framework, under clearly stated hypotheses, without need to introduce new *ad hoc* rules. Note also that the inference does not depend on multiplicative factors in the likelihood. This is one of the ways to state the likelihood principle, ideally desired by frequentists, but often violated. This 'principle' always and naturally holds in Bayesian statistics. It is important to remark that the use of unnecessary principles is dangerous, because there is a tendency to use them uncritically. For example, formulae resulting from maximum likelihood are often also used when non-uniform reasonable priors should be taken into account, or when the shape of  $\mathcal{L}(\theta; \mathbf{d})$  is far from being multi-variate Gaussian. (This is a kind of ancillary default hypothesis that comes together with this principle, and is the source of the often misused ' $\Delta(-\ln \mathcal{L}) = \frac{1}{2}$ ' rule to determine probability intervals.)

The usual least squares formulae are easily derived if we take the well-known case of pairs  $\{x_i, y_i\}$  (the generic  $\mathbf{d}$  stands for all data points) whose true values are related by a deterministic function  $\mu_{y_i} = y(\mu_{x_i}, \theta)$  and with Gaussian errors only in the ordinates; i.e. we

consider  $x_i \approx \mu_{x_i}$ . In the case of independence of the measurements, the likelihood-dominated result becomes

$$p(\boldsymbol{\theta} | \mathbf{x}, \mathbf{y}, I) \propto \prod_i \exp \left[ -\frac{(y_i - y(x_i, \boldsymbol{\theta}))^2}{2\sigma_i^2} \right] \quad (62)$$

or

$$p(\boldsymbol{\theta} | \mathbf{x}, \mathbf{y}, I) \propto \exp \left[ -\frac{1}{2}\chi^2 \right], \quad (63)$$

where

$$\chi^2(\boldsymbol{\theta}) = \sum_i \frac{(y_i - y(x_i, \boldsymbol{\theta}))^2}{\sigma_i^2} \quad (64)$$

is called 'chi-square', well known among physicists. Maximizing the likelihood is equivalent to minimizing  $\chi^2$ , and the most probable value of  $\boldsymbol{\theta}$  is easily obtained (i.e. the mode indicated with  $\boldsymbol{\theta}_m$ ), analytically in easy cases, or numerically for more complex ones.

As far as the uncertainty in  $\boldsymbol{\theta}$  is concerned, the widely used evaluation of the covariance matrix  $\mathbf{V}(\boldsymbol{\theta})$  (see section 5.6) from the Hessian,

$$(\mathbf{V}^{-1})_{ij}(\boldsymbol{\theta}) = \frac{1}{2} \frac{\partial^2 \chi^2}{\partial \theta_i \partial \theta_j} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_m}, \quad (65)$$

is merely a consequence of an assumed multi-variate Gaussian distribution of  $\boldsymbol{\theta}$ , that is, a parabolic shape of  $\chi^2$  (note that the ' $\Delta(-\ln \mathcal{L}) = \frac{1}{2}$ ' rule, and the ' $\Delta\chi^2 = 1$ ' rule', resulting from this has the same motivation). In fact, expanding  $\chi^2(\boldsymbol{\theta})$  in series around its minimum, we have

$$\chi^2(\boldsymbol{\theta}) \approx \chi^2(\boldsymbol{\theta}_m) + \frac{1}{2} \boldsymbol{\Delta}\boldsymbol{\theta}^T \mathbf{H} \boldsymbol{\Delta}\boldsymbol{\theta}, \quad (66)$$

where  $\boldsymbol{\Delta}\boldsymbol{\theta}$  stands for the the set of differences  $\theta_i - \theta_{m_i}$  and  $\mathbf{H}$  is the Hessian matrix, whose elements are given by twice the right-hand side of equation (65). Equation (63) then becomes

$$p(\boldsymbol{\theta} | \mathbf{x}, \mathbf{y}, I) \approx \propto \exp \left[ -\frac{1}{2} \boldsymbol{\Delta}\boldsymbol{\theta}^T \mathbf{H} \boldsymbol{\Delta}\boldsymbol{\theta} \right], \quad (67)$$

which we recognize to be a multi-variate Gaussian distribution if we identify  $\mathbf{H} = \mathbf{V}^{-1}$ . (The unusual symbol ' $\approx \propto$ ' in equation (67) stands for 'approximately proportional'.) After normalization, we finally get

$$p(\boldsymbol{\theta} | \mathbf{x}, \mathbf{y}, I) \approx (2\pi)^{-n/2} (\det \mathbf{V})^{-1/2} \exp \left[ -\frac{1}{2} \boldsymbol{\Delta}\boldsymbol{\theta}^T \mathbf{V}^{-1} \boldsymbol{\Delta}\boldsymbol{\theta} \right], \quad (68)$$

with  $n$  equal to the dimension of  $\boldsymbol{\theta}$  and  $\det \mathbf{V}$  indicating the determinant of  $\mathbf{V}$ . Holding this approximation,  $E(\boldsymbol{\theta})$  is approximately equal to  $\boldsymbol{\theta}_m$ . Note that the result (68) is exact when  $y(\mu_{x_i}, \boldsymbol{\theta})$  depends linearly on the various  $\theta_i$ .

In routine applications, the hypotheses that lead to the maximum likelihood and least squares formulae often hold. But when these hypotheses are not justified, we need to characterize the result by the multi-dimensional posterior distribution  $p(\boldsymbol{\theta})$ , going back to the more general expression equation (52).

The important conclusion from this section, as was the case for the definitions of probability in section 3, is that Bayesian methods often lead to well-known conventional results, but they are not introduced as new *ad hoc* rules. The analyst then acquires a heightened sense of awareness about the range of validity of the methods. One might as well use these 'recovered' methods within the Bayesian framework, with its more natural interpretation of the results. Then, one can speak about the uncertainty in the model parameters and quantify it with probability values, which is the usual way in which physicists think.

### 5.10. Gaussian approximation of the posterior distribution

The substance of the results seen in the previous section holds also in the case in which the prior is not flat and, hence, cannot be absorbed in the normalization constant of the posterior. In fact, in many practical cases the posterior exhibits an approximately (multi-variate) Gaussian shape, even if the prior was not trivial. Having at hand an un-normalized posterior  $\tilde{p}(\cdot)$ , i.e.

$$\tilde{p}(\boldsymbol{\theta} | \mathbf{d}, I) = p(\mathbf{d} | \boldsymbol{\theta}, I) p_0(\boldsymbol{\theta}, I), \quad (69)$$

we can take its minus-log function  $\varphi(\boldsymbol{\theta}) = -\ln \tilde{p}(\boldsymbol{\theta})$ . If  $\tilde{p}(\boldsymbol{\theta} | \mathbf{x}, \mathbf{y}, I)$  has an approximately Gaussian shape, it follows that

$$\varphi(\boldsymbol{\theta}) \approx \frac{1}{2} \boldsymbol{\Delta} \boldsymbol{\theta}^T \mathbf{V}^{-1} \boldsymbol{\Delta} \boldsymbol{\theta} + \text{constant}. \quad (70)$$

$\mathbf{V}$  can be evaluated as

$$(V^{-1})_{ij}(\boldsymbol{\theta}) \approx \left. \frac{\partial^2 \varphi}{\partial \theta_i \partial \theta_j} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_m}, \quad (71)$$

where  $\boldsymbol{\theta}_m$  was obtained from the minimum of  $\varphi(\boldsymbol{\theta})$ .

## 6. Uncertainties from systematic effects

The uncertainties described in the previous section are related to so-called random, or statistical errors. Other important sources are, generally speaking (see ISO 1993 for details), related to uncertain values of influence variables on which the observed values, or the data-analysis process, might depend. In physics, we usually refer to these as systematic effects or errors. They can be related to the parameters of the experiment, like a particle beam energy or the exposure time, or to environmental variables, like temperature and pressure, calibration constants of the detector, and all other parameters, ‘constants’ (in the physical sense), and hypotheses that enter the data analysis. The important thing is that we are unsure about their precise value. Let us denote all the influence variables by the vector  $\mathbf{h} = \{h_1, h_2, \dots, h_n\}$  and their joint pdf by  $p(\mathbf{h} | I)$ .

The treatment of uncertainties due to systematic errors has traditionally been lacking a consistent theory, essentially due to the unsuitability to standard statistical methods of dealing with uncertainty in the widest sense. Bayesian reasoning becomes crucial to handle these sources of uncertainty too, and even metrological organizations (ISO 1993) had to recognize it. For example, the ISO *type B* uncertainty is recommended to be ‘evaluated by scientific judgment based on all the available information on the possible variability’ (ISO 1993) of the influence quantities (see also D’Agostini (2003)).

### 6.1. Re-weighting of conditional inferences

The values of the influence variables and their uncertainties contribute to our background knowledge  $I$  about the experimental measurements. Using  $I_0$  to represent our very general background knowledge, the posterior pdf will then be  $p(\mu | \mathbf{d}, \mathbf{h}, I_0)$ , where the dependence on all possible values of  $\mathbf{h}$  has been made explicit. The inference that takes into account the uncertain vector  $\mathbf{h}$  is obtained using the rules of probability (see table 1) by integrating the joint probability over the uninteresting influence variables:

$$p(\mu | \mathbf{d}, I_0) = \int p(\mu, \mathbf{h} | \mathbf{d}, I_0) d\mathbf{h}, \quad (72)$$

$$p(\mu | \mathbf{d}, I_0) = \int p(\mu | \mathbf{d}, \mathbf{h}, I_0) p(\mathbf{h} | I_0) d\mathbf{h}. \quad (73)$$

As a simple, but important case, let us consider a single influence variable given by an additive instrumental offset  $z$ , which is expected to be zero because the instrument has been calibrated as well as is feasible and the remaining uncertainty is  $\sigma_z$ . Modelling our uncertainty in  $z$  as a Gaussian distribution with a standard deviation  $\sigma_z$ , the posterior for  $\mu$  is

$$p(\mu | d, I_0) = \int_{-\infty}^{+\infty} p(\mu | d, z, \sigma, I_0) p(z | \sigma_z, I_0) dz, \quad (74)$$

$$p(\mu | d, I_0) = \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi} \sigma} \exp\left[-\frac{(\mu - (d - z))^2}{2\sigma^2}\right] \frac{1}{\sqrt{2\pi} \sigma_z} \exp\left[-\frac{z^2}{2\sigma_z^2}\right] dz, \quad (75)$$

$$p(\mu | d, I_0) = \frac{1}{\sqrt{2\pi} \sqrt{\sigma^2 + \sigma_z^2}} \exp\left[-\frac{(\mu - d)^2}{2(\sigma^2 + \sigma_z^2)}\right]. \quad (76)$$

The result is that the net variance is the sum of the variances in the measurement and in the influence variable.

### 6.2. Joint inference and marginalization of nuisance parameters

A different approach, which produces identical results, is to think of a joint inference about both the quantities of interest and the influence variables:

$$p(\mu, \mathbf{h} | \mathbf{d}, I_0) \propto p(\mathbf{d} | \mu, \mathbf{h}, I_0) p_0(\mu, \mathbf{h} | I_0). \quad (77)$$

Then, marginalization is applied to the variables that we are not interested in (the so-called nuisance parameters), obtaining

$$p(\mu | \mathbf{d}, I_0) = \int p(\mu, \mathbf{h} | \mathbf{d}, I_0) d\mathbf{h}, \quad (78)$$

$$p(\mu | \mathbf{d}, I_0) \propto \int p(\mathbf{d} | \mu, \mathbf{h}, I_0) p_0(\mu, \mathbf{h} | I_0) d\mathbf{h}. \quad (79)$$

Equation (77) shows a peculiar feature of Bayesian inference, namely the possibility of making an inference about a number of variables larger than the number of the observed data. Certainly, there is no magic in it, and the resulting variables will be highly correlated. Moreover, the prior cannot be improper in all variables. But, by using informative priors in which experts feel confident, this feature allows one to tackle complex problems with missing or corrupted parameters. In the end, making use of marginalization, one can concentrate on the quantities of real interest.

The formulation of the problem in terms of equations (77) and (79) allows one to solve problems in which the influence variables might depend on the true value  $\mu$ , because  $p_0(\mu, \mathbf{h} | I_0)$  can model dependences between  $\mu$  and  $\mathbf{h}$ . In most applications,  $\mathbf{h}$  does not depend on  $\mu$ , and the prior factors into the product of  $p_0(\mu | I_0)$  and  $p_0(\mathbf{h} | I_0)$ . When this happens, we recover exactly the same results as those obtained using the re-weighting of conditional inferences approach described above.

### 6.3. Correlation in results caused by systematic errors

We can easily extend equations (73), (77), and (79) to a joint inference of several variables, which, as we have seen, are nothing but parameters  $\theta$  of suitable models. Using the alternative methods described in sections 6.1 and 6.2, we have

$$p(\theta | \mathbf{d}, \mathbf{h}, I_0) \propto p(\mathbf{d} | \theta, \mathbf{h}, I_0) p_0(\theta | I_0), \quad (80)$$

$$p(\theta | \mathbf{d}, I_0) = \int p(\theta | \mathbf{d}, \mathbf{h}, I_0) p(\mathbf{h} | I_0) d\mathbf{h}, \quad (81)$$

and

$$p(\boldsymbol{\theta}, \mathbf{h} | \mathbf{d}, I_0) \propto p(\mathbf{d} | \boldsymbol{\theta}, \mathbf{h}, I_0) p_0(\boldsymbol{\theta}, \mathbf{h} | I_0), \quad (82)$$

$$p(\boldsymbol{\theta} | \mathbf{d}, I_0) = \int p(\boldsymbol{\theta}, \mathbf{h} | \mathbf{d}, I_0) d\mathbf{h}, \quad (83)$$

respectively. The two methods lead to identical results, as can be seen by comparing equations (81) and (83).

Take a simple case of a common offset error of an instrument used to measure various quantities  $\mu_i$ , resulting in the measurements  $d_i$ . We model each measurement as  $\mu_i$  plus an error that is Gaussian distributed with a mean of zero and a standard deviation  $\sigma_i$ . The calculation of the posterior distribution can be performed analytically, with the following results (see D'Agostini (2003) for details):

- The uncertainty in each  $\mu_i$  is described by a Gaussian centred at  $d_i$ , with standard deviation  $\sigma(\mu_i) = \sqrt{\sigma_i^2 + \sigma_z^2}$ , consistent with equation (76).
- The joint posterior distribution  $p(\mu_1, \mu_2, \dots)$  does not factorize into the product of  $p(\mu_1)$ ,  $p(\mu_2)$ , etc, because correlations are automatically introduced by the formalism, consistent with the intuitive thinking of what a common systematic should do. Therefore, the joint distribution will be a multi-variate Gaussian that takes into account correlation terms.
- The correlation coefficient between any pair  $\{\mu_i, \mu_j\}$  is given by

$$\rho(\mu_i, \mu_j) = \frac{\sigma_z^2}{\sigma(\mu_i) \sigma(\mu_j)} = \frac{\sigma_z^2}{\sqrt{(\sigma_i^2 + \sigma_z^2)(\sigma_j^2 + \sigma_z^2)}}. \quad (84)$$

We see that  $\rho(\mu_i, \mu_j)$  has the behaviour expected from a common offset error; it is non-negative; it varies from practically zero, indicating negligible correlation, when ( $\sigma_z \ll \sigma_i$ ), to unity ( $\sigma_z \gg \sigma_i$ ), when the offset error dominates.

#### 6.4. Approximate methods and standard propagation applied to systematic errors

When we have many uncertain influence factors and/or the model of uncertainty is non-Gaussian, the analytic solution of equation (73), or equations (77)–(79) can be complicated, or must not exist at all. Then, numerical or approximate methods are needed. The most powerful numerical methods are based on MC techniques (see section 9 for a short account). This issue goes beyond the aim of this paper. In a recent comprehensive particle-physics paper by Ciuchini *et al* (2001), these ideas have been used to infer the fundamental parameters of the Standard Model of particle physics, using all available experimental information.

For routine use, a practical approximate method can be developed by thinking of the value inferred for the expected value of  $\mathbf{h}$  as a raw value, denoted by  $\mu_R$ , that is,  $\mu_R = \mu |_{\mathbf{h}=E(\mathbf{h})}$  ('raw' in the sense that it needs to be 'corrected' later for all possible values of  $\boldsymbol{\theta}$ , as will be clear in a while). The value of  $\mu$ , which depends on the possible values of  $\mathbf{h}$ , can be seen as a function of  $\mu_R$  and  $\mathbf{h}$ :

$$\mu = f(\mu_R, \mathbf{h}). \quad (85)$$

We have thus turned our inferential problem into a standard problem of evaluation of the pdf of a function of variables, of which the formulae to obtain approximate values for expectations and standard deviations in the case of independent input quantities (following the nomenclature of ISO 1993) are particularly known:

$$E(\mu) \approx f(E(\mu_R), E(\mathbf{h})), \quad (86)$$

$$\sigma^2(\mu) \approx \left( \frac{\partial f}{\partial \mu_R} \Big|_{E(\mu_R), E(h)} \right)^2 \sigma^2(\mu_R)q + \sum_i \left( \frac{\partial f}{\partial h_i} \Big|_{E(\mu_R), E(h)} \right)^2 \sigma^2(h_i). \quad (87)$$

Extension to multi-dimensional problems and treatment of correlations is straightforward (the well-known covariance matrix propagation) and we refer to D'Agostini and Raso (1999) for details. In particular, this reference contains approximate formulae valid up to second order, which allow us to take into account non-linearities relatively easily.

## 7. Comparison of models of different complexity

So far we have seen two typical inferential situations:

- (i) Comparison of simple models (section 4), where by simple we mean that the models do not depend on parameters to be tuned to the experimental data.
- (ii) Parametric inference given a model, to which we have devoted the last sections.

A more complex situation arises when we have several models, each of which might depend on several parameters. For simplicity, let us consider model  $A$  with  $n_A$  parameters  $\alpha$  and model  $B$  with  $n_B$  parameters  $\beta$ . In principle, the same Bayesian reasoning seen previously holds:

$$\frac{P(A | \text{Data}, I)}{P(B | \text{Data}, I)} = \frac{P(\text{Data} | A, I)}{P(\text{Data} | B, I)} \frac{P(A | I)}{P(B | I)}, \quad (88)$$

but we have to remember that the probability of the data, given a model, depends on the probability of the data, given a model and any particular set of parameters, weighted with the prior beliefs about parameters. We can use the same decomposition formula (see table 1) as already applied in treating systematic errors (section 6):

$$P(\text{Data} | M, I) = \int P(\text{Data} | M, \theta, I) p(\theta | I) d\theta, \quad (89)$$

with  $M = A, B$  and  $\theta = \alpha, \beta$ . In particular, the Bayes factor appearing in equation (88) becomes

$$\frac{P(\text{Data} | A, I)}{P(\text{Data} | B, I)} = \frac{\int P(\text{Data} | A, \alpha, I) p(\alpha | I) d\alpha}{\int P(\text{Data} | B, \beta, I) p(\beta | I) d\beta}, \quad (90)$$

$$\frac{P(\text{Data} | A, I)}{P(\text{Data} | B, I)} = \frac{\int \mathcal{L}_A(\alpha; \text{Data}) p_0(\alpha) d\alpha}{\int \mathcal{L}_B(\beta; \text{Data}) p_0(\beta) d\beta}. \quad (91)$$

The inference depends on the marginalized likelihood (89), also known as the evidence. Note that  $\mathcal{L}_M(\theta; \text{Data})$  has its largest value around the maximum likelihood point  $\theta_{\text{ML}}$ , but the evidence takes into account all prior possibilities of the parameters. Thus, it is not enough that the best fit of one model is superior to its alternative, in the sense that, for instance,

$$\mathcal{L}_A(\alpha_{\text{ML}}; \text{Data}) > \mathcal{L}_B(\beta_{\text{ML}}; \text{Data}), \quad (92)$$

and hence, assuming Gaussian models,

$$\chi_A^2(\alpha_{\min \chi^2}; \text{Data}) < \chi_B^2(\beta_{\min \chi^2}; \text{Data}), \quad (93)$$

to prefer model  $A$ . We have already seen that we need to take into account the prior beliefs in  $A$  and  $B$ . But even this is not enough: we also need to consider the space of possibilities and then the adaptation capability of each model. It is well understood that we do not choose an  $(n-1)$ -order polynomial as the best description—'best' in inferential terms—of  $n$  experimental points, though such a model always offers an exact pointwise fit. Similarly, we are much more impressed by, and we tend *a posteriori* to believe more in, a theory that absolutely predicts an

experimental observation, within a reasonable error, than another theory that performs similarly or even better after having adjusted many parameters.

This intuitive reasoning is expressed formally in equations (90) and (91). The evidence is given integrating the product  $\mathcal{L}(\theta)$  and  $p_0(\theta)$  over the parameter space. So, the more  $p_0(\theta)$  is concentrated around  $\theta_{ML}$ , the greater is the evidence in favour of that model. Instead, a model with a volume of parameter space much larger than the one selected by  $\mathcal{L}(\theta)$  gets disfavoured. The extreme limit is that of a hypothetical model with a sufficient number of parameters to describe whatever we shall observe. This effect is very welcome, and follows the Ockham's Razor scientific rule of discarding unnecessarily complicated models ('entities should not be multiplied unnecessarily'). This rule comes out of the Bayesian approach automatically and it is discussed, with examples of applications, in many papers. Berger and Jefferys (1992) introduce the connection between Ockham's Razor and Bayesian reasoning, and discuss the evidence provided by the motion of Mercury's perihelion in favour of Einstein's general relativity theory, compared with alternatives at that time. Examples of recent applications are Loredo and Lamb (2002) (analysis of neutrinos observed from supernova SN 1987A), John and Narlikar (2002) (comparisons of cosmological models), Hobson *et al* (2002) (combination of cosmological datasets), and Astone *et al* (2003) (analysis of coincidence data from gravitational wave detectors). These papers also give a concise account of underlying Bayesian ideas.

After having emphasized the merits of model comparison formalized in equations (90) and (91), it is important to mention a related problem. In parametric inference we have seen that we can make easy use of improper priors (see table 1), seen as limits of proper priors, essentially because they simplify in the Bayes formula. For example, we considered  $p_0(\mu | I)$  of equation (26) to be a constant, but this constant goes to zero as the range of  $\mu$  diverges. Therefore, it simplifies in equation (26), but not, in general, in equations (90) and (91), unless models  $A$  and  $B$  depend on the same number of parameters defined in the same ranges. Therefore, the general case of model comparison is limited to proper priors, and needs to be thought through better than when making parametric inference.

## 8. Choice of priors—a closer look

So far, we have considered mainly likelihood-dominated situations, in which the prior pdf can be included in the normalization constant. But one should be careful about the possibility of uncritically using uniform priors, as a 'prescription', or as a rule, though the rule might be associated with the name of famous persons. For instance, having made  $N$  interviews to infer the proportion  $\theta$  of a population that supports a party, it is not reasonable to assume a uniform prior of  $\theta$  between 0 and 1. Similarly, having to infer the rate  $r$  of a Poisson process (such that  $\lambda = rT$ , where  $T$  is the measuring time) related, for example, to proton decay, cosmic ray events, or gravitational wave signals, we do not believe, strictly, that  $p(r)$  is uniform between 0 and  $\infty$ . Besides natural physical cut-offs (e.g. very large proton decay  $r$  would prevent life, or even stars, to exist),  $p(r) = \text{constant}$  implies higher orders of magnitudes of  $r$  (see Astone and D'Agostini (1999) for details). In many cases (e.g. the mentioned searches for rare phenomena), our uncertainty could mean indifference over several orders of magnitude in the rate  $r$ . This indifference can be parametrized roughly with a prior uniform  $\ln r$  yielding  $p(r) \propto 1/r$  (the same prior is obtainable using invariant arguments, as will be shown later).

As the reader might imagine, the choice of priors is a much debated issue, also among Bayesians. We do not pretend to give definitive statements, but would just like to touch on some important issues concerning priors.

### 8.1. Logical and practical role of priors

Priors are pointed to by those critical of the Bayesian approach as the major weakness of the theory. Instead, Bayesians consider them a crucial and powerful key point of the method. Priors are logically crucial because they are necessary to make probability inversions via Bayes' theorem. This point remains valid even in the case in which they are vague and apparently disappear in the Bayes formula. Priors are powerful because they allow us to deal with realistic situations in which informative prior knowledge can be taken into account and properly balanced with the experimental information.

Indeed, we believe that one of the advantages of Bayesian analysis is that it explicitly admits the existence of prior information, which naturally leads to the expectation that the prior will be specified in any honest account of a Bayesian analysis. This crucial point is often obscured in other types of analysis, in large part because the analysts maintain that their method is 'objective'. Therefore, it is not easy, in those analyses, to recognize what are the specific assumptions made by the analyst—in practice the analyst's priors—and the assumptions included in the method (the latter assumptions are often unknown to the average practitioner).

### 8.2. Purely subjective assessment of prior probabilities

In principle, the point is simple, at least in the one-dimensional problem in which there is good perception of the possible range in which the uncertain variable of interest could lie: try your best to model your prior beliefs. In practice, this advice seems difficult to follow because, even if we have a rough idea of what the value of a quantity should be, the representation of the prior in mathematical terms seems very committal, because a pdf implicitly contains an infinite number of precise probabilistic statements. (Even the uniform distribution says that we believe exactly equally in all values. Who believes exactly that?) It is then important to understand that, when expressing priors, what matters is not the precise mathematical formula, but the gross value of the probability mass indicated by the formula, how probabilities are intuitively perceived and how priors influence posteriors. When we say, intuitively, we believe something with a 95% confidence, it means 'we are almost sure', but the precise value (95%, instead of 92% or 98%) is not very relevant. Similarly, when we say that the prior knowledge is modelled by a Gaussian distribution centred around  $\mu_0$  with standard deviation  $\sigma_0$  (equation (28)), it means that we are quite confident that  $\mu$  is within  $\pm\sigma_0$ , very sure that it is within  $\pm 2\sigma_0$  and almost certain that it is within  $\pm 3\sigma_0$ . Values even farther from  $\mu_0$  are possible, though we do not consider them very likely. But all models should be taken with a grain of salt, remembering that they are often just mathematical conveniences. For example, a textbook-Gaussian prior includes infinite deviations from the expected value and even negative values for physical quantities positively defined, like a temperature or a length—all absurdities, if taken literally. On the other hand, we think that all experienced physicists have in mind priors with low probability long tails in order to accommodate strong deviation from what is expected with highest probability. (Remember that where the prior is zero, the posterior must also be zero.)

Summing up this point, it is important to understand that a prior should indicate where the probability mass is concentrated, without taking too seriously the details, especially the tails of the distribution (which should however, be, extended enough to accommodate 'surprises'). The nice feature of Bayes' theorem is the ability to transform such vague, fuzzy priors into solid estimates, if a sufficient amount of good quality data are at hand. For this reason, the use of improper priors is not considered to be problematic. Indeed, improper priors can just be considered a convenient way of modelling relative beliefs.



In the case in which we have doubts about the choice of the prior, we can consider a family of functions with some hyperparameters. If we worry about the effect of the chosen prior on the posterior, we can perform a sensitivity analysis, i.e. we repeat the analysis for different, reasonable choices of the prior and check the variation of the result. The final uncertainty could, then, take into account also the uncertainty on the prior. Finally, in extreme cases in which priors play a crucial role and could dramatically change the conclusions, one should refrain from giving probabilistic results, providing, instead, only Bayes factors, or even just likelihoods. For an example of a recent result about gravitational wave searches presented in this way, see Astone *et al* (2002).

Having clarified the meaning and role of priors, it is rather evident that the practical choice of a prior depends on what is appropriate for the application. For example, in the area of imaging, smoothness of a reconstructed image might be appropriate in some situations. Smoothness may be imposed by a variety of means, for example, by simply setting the logarithm of the prior equal to an integral of the square of the second derivative of the image (von der Linden *et al* 1996b). A more sophisticated approach goes under the name of Markov random fields (MRF), which can even preserve sharp edges in the estimated images (Bouman and Sauer 1993, Saquib *et al* 1997). A similar kind of prior is often appropriate for deformable geometric models, which can be used to represent the boundaries between various regions, for example, organs in medical images (Cunningham *et al* 1998).

A procedure that helps in choosing the prior, especially important in the cases in which the parameters do not have a straightforwardly perceptible influence on data, is to build a prior predictive pdf and check whether this pdf would produce data conforming with our prior beliefs. The prior predictive distribution is the analogue of the (posterior) predictive distribution we encountered in section 5.7, with  $p(\boldsymbol{\theta} | \mathbf{d}, I)$  replaced by  $p(\boldsymbol{\theta} | I)$  (note that the example of section 5.7 was one-dimensional, with  $d_1 = d_f$  and  $\theta_1 = \mu$ ), i.e.  $p(\mathbf{d} | I) = \int p(\mathbf{d} | \boldsymbol{\theta}, I) p(\boldsymbol{\theta} | I) d\boldsymbol{\theta}$ .

Often, especially in complicated data analyses, we are not sufficiently knowledgeable about the details of the problem. Thus, informative priors have to be modelled that capture the judgement of experts. For example, Meyer and Booker (2001) show a formal process of prior elicitation which has the aim of reducing, as much as possible, the bias in the experts' estimates of their confidence limits. This approach allows us to combine the results from several experts. In short, we can suggest the use of the 'coherent bet' (section 2) to force experts to access their values of probability, asking them to provide an interval in which they feel 'practically sure', intervals on which they could wager 1 : 1, and so on.

### 8.3. Conjugate priors

Because of computational problems, modelling priors has traditionally been a compromise between a realistic assessment of beliefs and choosing a mathematical function that simplifies the analytical calculations. A well-known strategy is to choose a prior with a suitable form so that the posterior belongs to the same functional family as the prior. The choice of the family depends on the likelihood. A prior and posterior chosen in this way are said to be conjugate. For instance, given a Gaussian likelihood and choosing a Gaussian prior, the posterior is still Gaussian, as we have seen in equations (25), (28), and (29). This is because expressions of the form

$$K \exp \left[ -\frac{(x_1 - \mu)^2}{2\sigma_1^2} - \frac{(x_2 - \mu)^2}{2\sigma_2^2} \right]$$

can always be written in the form

$$K' \exp \left[ -\frac{(x' - \mu)^2}{2\sigma'^2} \right],$$

with suitable values for  $x'$ ,  $\sigma'$ , and  $K'$ . The Gaussian distribution is auto-conjugate. The mathematics is simplified but, unfortunately, only one shape is possible.

An interesting case, both for flexibility and practical interest is offered by the binomial likelihood (see section 5.3). Apart from the binomial coefficient,  $p(n | \theta, N)$  has the shape  $\theta^n(1 - \theta)^{N-n}$ , which has the same structure as the beta distribution, well known to statisticians:

$$\text{beta}(\theta | r, s) = \frac{1}{\beta(r, s)} \theta^{r-1} (1 - \theta)^{s-1}, \quad 0 \leq \theta \leq 1 \quad r, s > 0, \tag{94}$$

where  $\beta(r, s)$  stands for the beta function, defined as

$$\beta(r, s) = \int_0^1 \theta^{r-1} (1 - \theta)^{s-1} d\theta, \tag{95}$$

which can be expressed in terms of Euler's gamma function as  $\beta(r, s) = \Gamma(r) \Gamma(s) / \Gamma(r + s)$ . The expectation and variance of the beta distribution are

$$E(\theta) = \frac{r}{r + s}, \tag{96}$$

$$\sigma^2(\theta) = \frac{r s}{(r + s + 1)(r + s)^2} = E^2(\theta) \frac{s}{r} \frac{1}{r + s + 1}. \tag{97}$$

If  $r > 1$  and  $s > 1$ , then the mode is unique, and it is at  $\theta_m = (r - 1) / (r + s - 2)$ . Depending on the value of the parameters the beta pdf can take a large variety of shapes. For example, for large values of  $r$  and  $s$ , the function is very similar to a Gaussian distribution, while a constant function is obtained for  $r = s = 1$ . Using the beta pdf as the prior function in inferential problems with a binomial likelihood, we have

$$p(\theta | n, N, r, s) \propto [\theta^n (1 - \theta)^{N-n}] [\theta^{r-1} (1 - \theta)^{s-1}], \tag{98}$$

$$p(\theta | n, N, r, s) \propto \theta^{n+r-1} (1 - \theta)^{N-n+s-1}. \tag{99}$$

The posterior distribution is still a beta with  $r' = r + n$  and  $s' = s + N - n$ , and the expectation and standard deviation can be calculated easily from equations (96) and (97). These formulae demonstrate how the posterior estimates become progressively independent of the prior information in the limit of large numbers; this happens when both  $m \gg r$  and  $n - m \gg s$ . In this limit, we get the same result as for a uniform prior ( $r = s = 1$ ).

Table 2 lists some of the more useful conjugate priors. For a more complete collection of conjugate priors see, for example, Bernardo and Smith (1994) and Gelman *et al* (1995).

**Table 2.** Some useful conjugate priors.  $x$  and  $n$  stand for the observed value (continuous or discrete, respectively) and  $\theta$  is the generic symbol for the parameter to infer, corresponding to  $\mu$  of a Gaussian,  $\theta$  of a binomial and,  $\lambda$  of a Poisson distribution.

Likelihood	Conjugate prior	Posterior
$p(x   \theta)$	$p_0(\theta)$	$p(\theta   x)$
Normal( $\theta, \sigma$ )	Normal( $\mu_0, \sigma_0$ )	Normal( $\mu_1, \sigma_1$ ) (equations (30)–(32))
Binomial( $N, \theta$ )	Beta( $r, s$ )	Beta( $r + n, s + N - n$ )
Poisson( $\theta$ )	Gamma( $r, s$ )	Gamma( $r + n, s + 1$ )
Multinomial( $\theta_1, \dots, \theta_k$ )	Dirichlet( $\alpha_1, \dots, \alpha_k$ )	Dirichlet( $\alpha_1 + n_1, \dots, \alpha_k + n_k$ )

#### 8.4. General principle based priors

Many who advocate using the Bayesian approach still want to keep ‘subjective’ contributions to the inference to a minimum. Their aim is to derive prior functions based on ‘objective’ arguments or general principles. As the reader might guess, this subject is rather controversial, and the risk of transforming arguments, which might well be reasonable and useful in many circumstances, into dogmatic rules is high (D’Agostini 1999e).

*8.4.1. Transformation invariance.* An important class of priors arises from the requirement of transformation invariance. We shall consider two specific cases, translation invariance and scale invariance.

*Translation invariance.* Let us assume that we are indifferent about a transformation of the kind  $\theta' = \theta + b$ , where  $\theta$  is our variable of interest and  $b$  a constant. Then,  $p(\theta) d\theta$  is an infinitesimal mass element of probability for  $\theta$  to be in the interval  $d\theta$ . Translation invariance requires that this mass element remains unchanged when expressed in terms of  $\theta'$ , i.e.

$$p(\theta) d\theta = p(\theta') d\theta', \quad (100)$$

$$p(\theta) d\theta = p(\theta + b) d\theta, \quad (101)$$

since  $d\theta = d\theta'$ . It is easy to see that in order for equation (101) to hold for any  $b$ ,  $p(\theta)$  must be equal to a constant for all values of  $\theta$  from  $-\infty$  to  $+\infty$ . It is, therefore, an improper prior. As discussed above, this is just a convenient modelling. For practical purposes this prior should always be regarded as the limit for  $\Delta\theta \rightarrow \infty$  of  $p(\theta) = 1/\Delta\theta$ , where  $\Delta\theta$  is a large finite range around the values of interest.

*Scale invariance.* In other cases, we could be indifferent about a scale transformation, that is,  $\theta' = \beta\theta$ , where  $\beta$  is a constant. This invariance implies, since  $d\theta' = \beta d\theta$  in this case,

$$p(\theta) d\theta = p(\beta\theta)\beta d\theta, \quad (102)$$

i.e.

$$p(\beta\theta) = \frac{p(\theta)}{\beta}. \quad (103)$$

The solution of this functional equation is

$$p(\theta) \propto \frac{1}{\theta}, \quad (104)$$

as can be easily proved using equation (104) as the test solution in equation (103). This is the famous Jeffreys’ prior, since it was first proposed by Jeffreys. Note that this prior can also be stated as  $p(\log \theta) = \text{constant}$ , as can be easily verified. The requirement of scale invariance also produces an improper prior, in the range  $0 < \theta < \infty$ . Again, the improper prior must be understood as the limit of a proper prior extending several orders of magnitude around the values of interest. (Note that we constrain  $\theta$  to be positive because, traditionally, variables that are believed to satisfy this invariance are associated with positively defined quantities. Indeed, equation (104) has a symmetric solution for negative quantities.)

According to the supporters of these invariance motivated priors (see, e.g. Jaynes (1968), (1973), (1998), Sivia (1997), Fröhner (2000), Dose (2002)) variables associated with translation invariance are location parameters, such as the parameter  $\mu$  in a Gaussian model; variables associated with scale invariance are scale parameters, like the  $\sigma$  in a Gaussian model or  $\lambda$  in a Poisson model. For criticism about the (mis-)use of this kind of prior, see D’Agostini (1999d).

8.4.2. *Maximum-entropy priors.* Another principle-based approach to assigning priors is based on the Maximum Entropy principle (Jaynes 1957a, 1983, 1998, Tribus 1969, von der Linden 1995, Sivia 1997, Fröhner 2000). The basic idea is to choose the prior function that maximizes the Shannon–Jaynes information entropy:

$$S = - \sum_i^n p_i \ln p_i, \quad (105)$$

subject to whatever is assumed to be known about the distribution. The larger  $S$  is, the greater is our ignorance about the uncertain value of interest. The value  $S = 0$  is obtained for a distribution that concentrates all the probability into a single value. In the case of no constraint other than normalization ( $\sum_i^n p_i = 1$ ),  $S$  is maximized by the uniform distribution,  $p_i = 1/n$ , which is easily proved using Lagrange multipliers. For example, if the variable is an integer between 0 and 10, a uniform distribution  $p(x_i) = \frac{1}{11}$  gives  $S = 2.40$ . Any binomial distribution with  $n = 10$  gives a smaller value, with a maximum of  $S = 1.88$  for  $\theta = \frac{1}{2}$  and a limit of  $S = 0$  for  $\theta \rightarrow 0$  or 1, where  $\theta$  is now the parameter of the binomial that gives the probability of success at each trial.

Two famous cases of Maximum Entropy priors for continuous variables are those relating to the cases when the only information about the distribution is either the expected value or the expected value and the variance. Indeed, these are special cases of general constraints on the moments of distribution  $M_r$  (see table 1). For  $r = 0$  and 1,  $M_r$  is equal to unity and to the expected value, respectively. The first and second moments together provide the variance (see table 1 and section 5.6). Let us sum up what the assumed knowledge of the various moments provides (see, e.g. Sivia (1997), Dose (2002)).

*Knowledge about  $M_0$ .* Normalization alone provides a uniform distribution over the interval in which the variable is defined:

$$p(\theta | M_0 = 1) = \frac{1}{b - a} \quad a \leq \theta \leq b. \quad (106)$$

This is the extension to continuous variables of the discrete case we saw earlier.

*Knowledge about  $M_0$  and  $M_1$  (i.e. about  $E(\theta)$ ).* Adding to the constraint  $M_0 = 1$  the knowledge about the expectation of the variable, plus the requirement that all non-negative values are allowed, an exponential distribution is obtained:

$$p(\theta | M_0 = 1, M_1 \equiv E(\theta)) = \frac{1}{E(\theta)} e^{-\theta/E(\theta)} \quad 0 \leq \theta < \infty. \quad (107)$$

*Knowledge about  $M_0$ ,  $M_1$ , and  $M_2$  (i.e. about  $E(\theta)$  and  $\sigma(\theta)$ ).* Finally, the further constraint provided by the standard deviation (related to the first and second moments by the equation  $\sigma^2 = M_2 - M_1^2$ ) yields a prior with a Gaussian shape independently of the range of  $\theta$ , i.e.

$$p(\theta | M_0 = 1, E(\theta), \sigma(\theta)) = \frac{\exp[-(\theta - E(\theta))^2/2\sigma^2(\theta)]}{\int_a^b \exp[-(\theta - E(\theta))^2/2\sigma^2(\theta)] d\theta} \quad a \leq \theta \leq b. \quad (108)$$

The standard Gaussian is recovered when  $\theta$  is allowed to be any real number.

Note, however, that the counterpart of equation (105) for continuous variables is not trivial, since all  $p_i$  of equation (105) tend to zero. Hence the analogous functional form  $\int p(\theta) \ln p(\theta) d\theta$  no longer has a sensible interpretation in terms of uncertainty, as remarked by Bernardo and Smith (1994). The Jaynes' solution is to introduce a 'reference' density  $m(\theta)$  to make entropy invariant under coordinate transformation via  $\int p(\theta) \ln[p(\theta)/m(\theta)] d\theta$ .

(It is important to note that the first and the third cases discussed above are valid under the assumption of a unity reference density.) This solution is not universally accepted (see Bernardo and Smith (1994)), even though it conforms to the requirements of dimensional analysis. Anyhow, besides formal aspects and the undeniable aid of Maximum Entropy methods in complicated problems such as image reconstruction (Buck and Macaulay 1991), we find it very difficult, if not impossible, to believe that a practitioner holds that status of knowledge which gives rise to the two celebrated cases discussed earlier. We find the approach described in section 8.2, which goes the other way around, more reasonable; we have a rough idea of where the quantity of interest could be, then we try to model it and to summarize it in terms of the expected value and standard deviation. In particular, we find untenable the position of those who state that the Gaussian distribution can only be justified by the Maximum Entropy principle.

*8.4.3. Reference priors.* We conclude this discussion on priors by mentioning ‘reference analysis’, which is an area of active research among statisticians. The intention is, similarly to that for other priors motivated by basic principles, that of ‘characterizing a “non-informative” or “objective” prior distribution, representing “prior ignorance”, “vague prior knowledge”, and “letting the data speak for themselves”’ (Bernardo and Smith 1994). However, ‘the problem is more complex than the apparent intuitive immediacy of these words and phrases would suggest’ (Bernardo and Smith (1994), p 298):

Put bluntly: data cannot ever speak entirely for themselves: every prior specification has some informative posterior or predictive implications; and ‘vague’ is itself much too vague an idea to be useful. There is no ‘objective’ prior that represents ignorance.

On the other hand, we recognize that there is often a pragmatically important need for a form of prior to posterior analysis capturing, in some well-defined sense, the notion of the prior having a minimal effect, relative to the data, on the final inference. Such a reference analysis might be required as an approximation to actual beliefs; more typically, it might be required as a limiting ‘what if?’ baseline in considering a range of prior to posterior analyses, or as a default option when there are insufficient resources for detailed elicitation of actual prior knowledge.

... From the approach we adopt, it will be clear that the reference prior component of the analysis is simply a mathematical tool. It has considerable pragmatic importance in implementing a reference analysis, whose role and character will be precisely defined, but it is not a privileged, ‘unique non-informative’ or ‘objective’ prior.

The curious reader may take a look at Bernardo and Smith (1994) and references therein, as well as at Bernardo (1997).

## 9. Computational issues

The application of Bayesian ideas leads to computational problems, mostly related to the calculation of integrals for normalizing the posterior pdf and for obtaining credibility regions, or simply the moments of distribution (and, hence, expectations, variances, and covariances). The difficulties become challenging for problems involving many parameters. This is one of the reasons why Bayesian inference was abandoned at the beginning of the twentieth century in favour of simplified—and simplistic—methods. Indeed, the Bayesian renaissance over the past few decades is largely due to the emergence of new numerical methods and the dramatic

increases in computational power, along with work on the foundations of the theory, which clarified various issues.

### 9.1. Overview of approximate computational strategies

In earlier sections we have already seen some ‘tricks’ for simplifying the calculations. The main topic of this section will be an introduction to MC. But, before doing that, we think it is important to summarize the various ‘tricks’ here. Much specialized literature is available on several aspects of computation in statistics (for an excellent review paper on the subject see Smith (1991)).

*Conjugate priors.* We discussed this topic in section 8.3, giving a couple of typical simple examples and references for a more detailed list of famous conjugate distributions. We would like to remark here that a conjugate prior is a special case of the class of priors that simplify the calculation of the posterior (the uniform prior is the simplest of this kind of prior).

*Gaussian approximation.* For reasons that are connected with the central limit theorem, when there is a large amount of consistent data, the posterior tends to be Gaussian, practically independently of the exact shape of the prior. The (multi-variate) Gaussian approximation, which we encountered in section 5.10, has an important role for applications, either as a reasonable approximation of the ‘true’ posterior, or as a starting point for searching for a more accurate description of it. We also saw that in the case of practically flat priors this method recovers the well-known minimum chi-square or maximum likelihood methods.

*Numerical integration.* In the case of low-dimensional problems, standard numerical integration using either scientific library functions or the interactive tools of modern computer packages provide an easy solution to many problems (mention must also be made of the graphical capabilities of modern programs, which allow the shape of the posterior to be inspected and the best calculation strategy decided upon). This is a vast and growing subject, into which we cannot enter in any depth here, but we assume the reader is familiar with some of these programs or packages.

*MC methods.* MC methodology is a science in itself and it is way beyond our remit to provide an exhaustive introduction to it here. Nevertheless, we would like to introduce briefly some ‘modern’ (though the seminal work is already half a century old) methods that are becoming extremely popular and are often associated with Bayesian analysis, the so-called Markov Chain MC (MCMC) methods.

### 9.2. MC methods

*9.2.1. Estimating expectations by sampling.* The easy part of the Bayesian approach is to write down the un-normalized distribution of the parameters (section 5.10), given the prior and the likelihood. This is simply  $\tilde{p}(\boldsymbol{\theta} | \mathbf{d}, I) = p(\mathbf{d} | \boldsymbol{\theta}, I) p(\boldsymbol{\theta} | I)$ . The difficult task is to normalize this function and to calculate all expectations in which we are interested, such as expected values, variances, covariances, and other moments. We might also want to get marginal distributions, credibility intervals (or hypervolumes) and so on. As is well known, if we were able to sample the posterior (even the un-normalized one), i.e. to generate points of the parameter space according to their probability, we would have solved the problem, at least approximately. For example, the one-dimensional histogram of parameter  $\theta_i$  would represent its marginal and would allow the calculation of  $E(\theta_i) \approx \langle \theta_i \rangle$ ,  $\sigma(\theta_i) \approx \langle \theta_i^2 \rangle - \langle \theta_i \rangle^2$  and of probability intervals ( $\langle \cdot \rangle$  in the previous formulae stand for arithmetic averages of the MC sample).

Let us consider the probability function  $p(\mathbf{x})$  of the discrete variables  $\mathbf{x}$  and a function  $f(\mathbf{x})$  for which we want to evaluate the expectation over the distribution  $p(\mathbf{x})$ . Extending the one-dimensional formula of table 1 to  $n$  dimensions we have

$$E[f(\mathbf{x})] = \sum_{x_1} \cdots \sum_{x_n} f(x_1, \dots, x_n) p(x_1, \dots, x_n), \quad (109)$$

$$E[f(\mathbf{x})] = \sum_i f(\mathbf{x}_i) p(\mathbf{x}_i), \quad (110)$$

where the summation in equation (109) is over the components, while the summation in equation (110) is over possible points in the  $n$ -dimensional space of the variables. The result is the same.

If we are able to sample a large number of points  $N$  according to the probability function  $p(\mathbf{x})$ , we expect each point to be generated  $m_i$  times. The average  $\langle f(\mathbf{x}) \rangle$ , calculated from the sample as

$$\langle f(\mathbf{x}) \rangle = \frac{1}{N} \sum_t f(\mathbf{x}_t), \quad (111)$$

(in which  $t$  is used for the index as a reminder that this is a sum over a ‘time’ sequence) can also be rewritten as

$$\langle f(\mathbf{x}) \rangle = \sum_i f(\mathbf{x}_i) \frac{m_i}{N}, \quad (112)$$

just grouping together the outcomes giving the same  $\mathbf{x}_i$ . For a very large  $N$ , the ratios  $m_i/N$  are expected to be ‘very close’ to  $p(\mathbf{x}_i)$  (Bernoulli’s theorem), and thus  $\langle f(\mathbf{x}) \rangle$  becomes a good approximation of  $E[f(\mathbf{x})]$ . In fact, this approximation can be good (within tolerable errors) even if not all  $m_i$  are large and, indeed, even if many of them are null. Moreover, the same procedure can be extended to the continuum, in which case ‘all points’ ( $\infty^n$ ) can never be sampled.

For simple distributions, there are well-known, standard techniques for generating pseudo-random numbers starting from pseudo-random numbers distributed uniformly between 0 and 1 (computer libraries are available for sampling points according to the most common distributions). We shall not enter into these basic techniques, but will concentrate instead on the calculation of expectations in more complicated cases.

**9.2.2. Rejection sampling.** Let us assume that we are able to generate points according to some function  $g(\mathbf{x})$ , such that, given a constant  $c$ ,  $p(\mathbf{x}) \leq cg(\mathbf{x})$ . We generate  $\mathbf{x}^*$  according to  $g(\mathbf{x})$  and decide to accept it with probability  $p(\mathbf{x}^*)/cg(\mathbf{x}^*)$  (i.e. we extract another random number between 0 and 1 and accept the point if this number is below that ratio). It is easy to show that this procedure reshapes  $g(\mathbf{x})$  to  $p(\mathbf{x})$  and that it does not depend on the absolute normalization of  $p(\mathbf{x})$  (any normalization constant can be absorbed in the multiplicative constant  $c$ ). A trivial choice of  $g(\mathbf{x})$ , especially for simple one-dimensional problems, is a uniform distribution (this variation is known as the hit or miss method), though clearly it can be very inefficient.

**9.2.3. Importance sampling.** In this method, too, we start from an ‘easy’ function  $g(\mathbf{x})$ , which ‘we hope’ will approximate the  $p(\mathbf{x})$  of interest, of which in fact we know only its un-normalized expression  $\tilde{p}(\mathbf{x})$ . However, there is no requirement about how  $g(\mathbf{x})$  approximates  $p(\mathbf{x})$  (but the goodness of approximation will have an impact on the efficacy of the method), apart from the condition that  $g(\mathbf{x}_i)$  must be positive wherever  $p(\mathbf{x}_i)$  is positive.

The function  $g(\mathbf{x})$  can be used in the calculation of  $E[f(\mathbf{x})]$ , if we notice that  $E[f(\mathbf{x})]$  can be rewritten as follows:

$$E[f(\mathbf{x})] = \frac{\int f(\mathbf{x}) \tilde{p}(\mathbf{x}) \, d\mathbf{x}}{\int \tilde{p}(\mathbf{x}) \, d\mathbf{x}}, \quad (113)$$

$$E[f(\mathbf{x})] = \frac{\int f(\mathbf{x}) [\tilde{p}(\mathbf{x})/g(\mathbf{x})] g(\mathbf{x}) \, d\mathbf{x}}{\int [\tilde{p}(\mathbf{x})/g(\mathbf{x})] g(\mathbf{x}) \, d\mathbf{x}}, \quad (114)$$

$$E[f(\mathbf{x})] = \frac{E_g[f(\mathbf{x}) \tilde{p}(\mathbf{x})/g(\mathbf{x})]}{E_g[\tilde{p}(\mathbf{x})/g(\mathbf{x})]}, \quad (115)$$

where the the symbol  $E_g$  is a reminder that the expectation is calculated over the distribution  $g(\mathbf{x})$ . Finally, the strategy can be implemented in MC using equation (111) for the two expectations:

$$E[f(\mathbf{x})] \approx \frac{\sum_t f(\mathbf{x}_t) \tilde{p}(\mathbf{x}_t)/g(\mathbf{x}_t)}{\sum_t \tilde{p}(\mathbf{x}_t)/g(\mathbf{x}_t)}. \quad (116)$$

From the same sample it is also possible to evaluate the normalization constant, given by the denominator of equation (113), i.e.

$$Z = \int \tilde{p}(\mathbf{x}) \, d\mathbf{x} \approx \frac{1}{N} \sum_t \frac{\tilde{p}(\mathbf{x}_t)}{g(\mathbf{x}_t)}. \quad (117)$$

The computation of this quantity is particularly important when we are dealing with model comparison and  $Z$  has the meaning of ‘evidence’ (section 7).

It is easy to see that the method works well if  $g(\mathbf{x})$  overlaps well with  $p(\mathbf{x})$ . Thus, a proper choice of  $g(\mathbf{x})$  can be made by studying where the probability mass of  $p(\mathbf{x})$  is concentrated (e.g. finding the mode of the distribution in a numerical way). Often a Gaussian function is used for  $g(\mathbf{x})$ , with parameters chosen to approximate  $p(\mathbf{x})$  in the proximity of the mode, as described in section 5.10. In other cases, other functions can be used which have more pronounced tails, like  $t$ -Student or Cauchy distributions. Special techniques, into which we cannot enter here, allow  $n$  independent random numbers to be generated and, subsequently, by proper rotations, turned into other numbers which have a correlation matrix equal to that of the multi-dimensional Gaussian that approximates  $p(\mathbf{x})$ .

Note, finally, that, contrary to the rejection sampling, importance sampling is not suitable for generating samples of ‘unweighted events’, such as those routinely used in the planning and analysis of many kinds of experiment, especially particle physics experiments.

**9.2.4. Metropolis algorithm.** A different class of MC methods is based on Markov chains and is known as MCMC. The basic difference from the methods described above is that the sequence of generated points takes a kind of random walk in parameter space, instead of each point being generated, one independent of another. Moreover, the probability of jumping from one point to another depends only on the last point and not on the entire previous history (this is the peculiar property of a Markov chain). There are several MCMC algorithms. One of the most popular and simplest algorithms, applicable to a wide class of problems, is the Metropolis algorithm (Metropolis *et al* 1953). One starts from an arbitrary point  $\mathbf{x}_0$  and generates the sequence by repeating the following cycle, with  $\mathbf{x}_t$  being the previously selected point at each iteration:

- (i) Select a new trial point  $\mathbf{x}^*$ , chosen according to a symmetric proposal pdf  $q(\theta^* | \theta_t)$ .



(ii) Calculate the acceptance probability

$$A(\mathbf{x}^* | \mathbf{x}_t) = \min \left[ 1, \frac{\tilde{p}(\boldsymbol{\theta}^*)}{\tilde{p}(\boldsymbol{\theta}_t)} \right]. \quad (118)$$

(iii) Accept  $\mathbf{x}^*$  with probability  $A(\mathbf{x}_t, \mathbf{x}^*)$ , i.e.

- if  $\tilde{p}(\boldsymbol{\theta}^*) \geq \tilde{p}(\boldsymbol{\theta}_t)$ , then accept  $\mathbf{x}^*$ ;
- if  $\tilde{p}(\boldsymbol{\theta}^*) < \tilde{p}(\boldsymbol{\theta}_t)$ , extract a uniform random number between 0 and 1 and accept  $\mathbf{x}^*$  if the random number is less than  $\tilde{p}(\boldsymbol{\theta}^*)/\tilde{p}(\boldsymbol{\theta}_t)$ .

If the point is accepted, then  $\mathbf{x}_{t+1} = \mathbf{x}^*$ . Otherwise  $\mathbf{x}_{t+1} = \mathbf{x}_t$ .

This algorithm allows a jump  $\mathbf{x}_t$  to  $\mathbf{x}_{t+1}$  with probability  $T(\mathbf{x}_{t+1} | \mathbf{x}_t)$  (the transition kernel) equal to  $A(\mathbf{x}^* | \mathbf{x}_t) q(\boldsymbol{\theta}^* | \boldsymbol{\theta}_t)$ . The algorithm has a stationary asymptotic distribution equal to  $p(\mathbf{x})$  (i.e. the chain is ergodic) and ensures detailed balance:

$$p(\mathbf{x}_{t+1}) T(\mathbf{x}_t | \mathbf{x}_{t+1}) = p(\mathbf{x}_t) T(\mathbf{x}_{t+1} | \mathbf{x}_t). \quad (119)$$

By construction, the algorithm does not depend on the normalization constant, since what matters is the ratio of the pdfs. The variation of the algorithm in which the proposal pdf  $q(\cdot)$  is not symmetric is due to Hastings (1970) and for this reason the algorithm is often also called Metropolis–Hasting. Moreover, what has been described here is the global Metropolis algorithm, in contrast to the local one, in which a cycle affects only one component of  $\mathbf{x}$ .

The fact that this algorithm belongs to the class of MCMC gives rise to two problems. First, each point in the chain has some correlation with the points that immediately preceded it, and usually the chain moves slowly (and irregularly) from one region in the variable space to another (note also that, if a proposed point is not accepted, the chain keeps the same position in the next step, and this circumstance can happen several times consecutively). Second, the initial part of the sequence is strongly influenced by the arbitrary starting point. Therefore, it is necessary to remove the initial part of the chain.

Using an MCMC for a complex problem is not an automatic procedure and some tuning is needed. One of the important things to choose with care is the proposal function. If the jumps proposed are too small the chain moves too slowly and can even remain trapped in a subregion and never sample the rest of the parameter space, if the probability distribution is defined over disconnected regions. If too large steps are proposed, the proposed points could often fall in very low probability regions and not be accepted, in which case the chain remains stuck at a point for many cycles.

For an interesting, insightful introduction to the principles and applications of MCMC see Kass *et al* (1998). A nice tutorial is given by Hanson (2000). A recent application of Bayesian methods in cosmology, which uses MCMC and contains a pedagogical introduction too, can be found in Lewis and Bridle (2002). For a good treatise, freely available on the web, Neel (1993) is recommended. The reader will find that MCMC techniques are used to solve complex problems graphically represented in terms of Bayesian networks (also known as belief networks or simply probabilistic networks). This subject, which has revolutionized the way of thinking about artificial intelligence and the uncertainty issues related to it, goes beyond the purpose of this paper. The interested reader can find more information in Pearl (1988), BUGS (1996), Cowell *et al* (1999), Cozman (2001) and references therein.

## 10. Conclusions

The gain in popularity that Bayesian methods have enjoyed in recent years is due to various conceptual and practical advantages they have over other approaches, among which are the

following:

- the recovery of the intuitive idea of probability as a valid concept for treating scientific problems;
- the simplicity and naturalness of the basic tool;
- the capability of combining prior knowledge and experimental information;
- the property permitting automatic updating as soon as new information becomes available;
- the transparency of the methods, which allow the different assumptions upon which an inference may depend to be checked and changed;
- the high degree of awareness the methods give to the user.

In this paper, we have seen how to build a theory of uncertainty in measurement as a straightforward application of the basic Bayesian ideas, without unnecessary principles or *ad hoc* prescriptions. In particular, the uncertainty due to systematic errors can be treated in a consistent and powerful way.

Providing an exact solution for inferential problems can easily lead to computational difficulties. We have seen several methods by which such difficulties can be overcome, either by using suitable approximations, or by using modern computational methods. In particular, it has been shown that the approximate solution often coincides with a 'conventional' method, but only under well-defined conditions. Thus, for example, minimum  $\chi^2$  formulae can be used, with a Bayesian spirit and with a natural interpretation of the results, in all those routine cases in which the analyst considers as reasonable the conditions of their validity.

A variety of examples of applications have been shown, or mentioned, in this paper. Nevertheless, the aim of the author was not to provide a complete review of Bayesian methods and applications, but rather to introduce those Bayesian ideas that could be of help in understanding more specialized literature. Compendia of the Bayesian theory are given in Bernardo and Smith (1994), O'Hagan (1994), Robert (2001). Classic, influential books are Jeffreys (1961), de Finetti (1974), Jaynes (1998). Among the many books introducing Bayesian methods, Sivia (1996) is particularly suitable for physicists. Other recommended texts which treat general aspects of data analysis are Box and Tiao (1973), Bretthorst (1988), Lee (1989), Gelman *et al* (1995), Cowell *et al* (1999), Denison *et al* (2002), Press (2002). More specific applications can be found in the proceedings of the conference series and several web sites. Some useful starting points for web navigation are given:

ISBA book list	<a href="http://www.bayesian.org/books/books.html">http://www.bayesian.org/books/books.html</a>
UAI proceedings	<a href="http://www2.sis.pitt.edu/dsl/UAI/uai.html">http://www2.sis.pitt.edu/dsl/UAI/uai.html</a>
BIPS	<a href="http://astrosun.tn.cornell.edu/staff/loredo/bayes/">http://astrosun.tn.cornell.edu/staff/loredo/bayes/</a>
BLIP	<a href="http://www.ar-tiste.com/blip.html">http://www.ar-tiste.com/blip.html</a>
IPP Bayesian analysis group	<a href="http://www.ipp.mpg.de/OP/Datenanalyse/">http://www.ipp.mpg.de/OP/Datenanalyse/</a>
Valencia meetings	<a href="http://www.uv.es/~bernardo/valenciam.html">http://www.uv.es/~bernardo/valenciam.html</a>
Maximum Entropy resources	<a href="http://omega.albany.edu:8008/maxent.html">http://omega.albany.edu:8008/maxent.html</a>
MCMC preprint service	<a href="http://www.statslab.cam.ac.uk/~mcmc/">http://www.statslab.cam.ac.uk/~mcmc/</a>

### Acknowledgments

I am indebted to Volker Dose and Ken Hanson for extensive discussions concerning the contents of this paper, as well as for substantial editorial help. The manuscript has also benefited from comments by Tom Lored.

## References

- Astone P *et al* 2002 Search for correlation between GRB's detected by BeppoSAX and gravitational wave detectors EXPLORER and NAUTILUS *Phys. Rev.* **66** 102002
- Astone P and D'Agostini G 1999 Inferring the intensity of Poisson processes at limit of detector sensitivity (with a case study on gravitational wave burst search) *CERN-EP/99-126*
- Astone P, D'Agostini G and D'Antonio S 2003 Bayesian model comparison applied to the Explorer-Nautilus 2001 coincidence data, arXiv:gr-qc/0304096
- Babu G J and Feigelson E D (ed) 1992 *Statistical Challenges in Modern Astronomy I* (New York: Springer)
- Babu G J and Feigelson E D (ed) 1997 *Statistical Challenges in Modern Astronomy II* (New York: Springer)
- Berger J O and Jefferys W H 1992 Sharpening Ockham's razor on a Bayesian strop *Am. Sci.* **89** 64–72
- 1992 *J. Ital. Stat. Soc.* **1** 17
- Bernardo J M, Berger J O, Dawid A P and Smith A F M (ed) 1999 *Bayesian Statistics 6* (Oxford: Oxford University)
- Bernardo J M and Smith F M 1994 *Bayesian Theory* (Chichester: Wiley)
- Bernardo J M 1997 Non-informative priors do not exist *J. Stat. Planning Infer.* **65** 159
- Bontekoe T R, Koper E and Kester D J M 1994 Pyramid maximum entropy images of IRAS survey data *Astron. Astrophys.* **284** 1037
- Bouman C A and Sauer K 1993 A generalized Gaussian image model for edge-preserving MAP estimation *IEEE Trans. Image Process.* **2** 296
- Box G E P and Tiao G C 1973 *Bayesian Inference in Statistical Analysis* (Chichester: Wiley)
- Bretthorst G L 1988 *Bayesian Spectrum Analysis and Parameter Estimation* (Berlin: Springer)
- BUGS 1996 <http://www.mrc-bsu.cam.ac.uk/bugs/welcome.shtml>
- Buck B and Macaulay V A (ed) 1991 *Maximum Entropy in Action* (Oxford: Oxford University Press)
- Ciuchini M *et al* 2001 2000 CKM-triangle analysis: a critical review with updated experimental inputs and theoretical parameters *J. High Energy Phys.* **0107** 013
- Coletti G and Scozzafava R 2002 *Probabilistic Logic in a Coherent Setting* (Dordrecht: Kluwer)
- Cousins R D 1995 Why isn't every physicist a Bayesian? *Am. J. Phys.* **63** 398
- Cowell R G, Dawid A P, Lauritzen S L and Spiegelhalter D J 1999 *Probabilistic Networks and Expert Systems* (New York: Springer)
- Cox R T 1946 Probability, frequency and reasonable expectation *Am. J. Phys.* **14** 1
- Cozman F B 2001 JavaBayes version 0.346—Bayesian networks in Java <http://www-2.cs.cmu.edu/~javabayes/Home/>
- Cunningham G S, Hanson K M and Battle X L 1998 Three-dimensional reconstructions from low-count SPECT data using deformable models *Opt. Expr.* **2** 227
- D'Agostini G 1999a Bayesian reasoning versus conventional statistics in high energy physics *Maximum Entropy and Bayesian Methods* ed von der Linden W *et al* (Dordrecht: Kluwer)
- D'Agostini G 1999b Sceptical combination of experimental results: General considerations and application to epsilon-prime/epsilon *CERN-EP/99-139*
- D'Agostini G 1999c Bayesian reasoning in high-energy physics: principles and applications *CERN Report* 99–03
- D'Agostini G 1999d Teaching statistics in the physics curriculum: Unifying and clarifying role of subjective probability *Am. J. Phys.* **67** 1260
- D'Agostini G 1999e Overcoming prior anxiety *Bayesian Methods in the Sciences* ed J M Bernardo; special issue of *Rev. Acad. Cien. (Madrid)* **93** 311
- D'Agostini G 2000 Confidence limits: what is the problem? Is there the solution? *CERN Report* 2000-005, ed F James and L Lyons (Geneva: CERN) p 3
- D'Agostini G 2002 Minimum bias legacy of search results 2002 *Nucl. Phys. Proc. Suppl.* **109** 148
- D'Agostini G 2003 *Bayesian Reasoning in Data Analysis—A Critical Introduction* (Singapore: World Scientific) (reviewed and extended version of D'Agostini 1999c)
- D'Agostini G and Degrossi G 1999 Constraints on the Higgs boson mass from direct searches and precision measurements *Eur. Phys. J.* **C10** 663
- D'Agostini G and Raso M 2000 Uncertainties due to imperfect knowledge of systematic effects: general considerations and approximate formulae *CERN-EP/2000-026*
- de Finetti B 1974 *Theory of Probability* (Chichester: Wiley)
- Denison D G T, Holmes C C, Mallick B K and Smith A F M 2002 *Bayesian Methods for Nonlinear Classification and Regression* (Chichester: Wiley)
- DIN (Deutsches Institut für Normung) 1996 *Grundlagen der Messtechnik—Teil 3: Auswertung von Messungen einer einzelnen Messgröße, Messunsicherheit* DIN 1319-3 (Berlin: Beuth Verlag)

- DIN (Deutsches Institut für Normung) 1999 Grundlagen der Messtechnik—Teil 4: Auswertung von Messungen, Messunsicherheit DIN 1319-4 (Berlin: Beuth Verlag)
- Dose V 2002 Bayes in five days, CIPS, MPI für Plasmaphysik, Garching, Germany, Reprint 83, May 2002
- Dose V and von der Linden W 1999 Outlier tolerant parameter estimation *Maximum Entropy and Bayesian Methods* ed von der Linden W *et al* (Dordrecht: Kluwer) p 47
- Efron B 1986a Why isn't everyone a Bayesian? *Am. Stat.* **40** 1
- Efron B 1986b Reply to Zellner 1986 *Am. Stat.* **40** 331
- Fischer R, Mayer M, von der Linden W and Dose V 1997 Enhancement of the energy resolution in ion-beam experiments with the maximum-entropy method *Phys. Rev. E* **55** 6667
- Fischer R, Mayer M, von der Linden W and Dose V 1998 Energy resolution enhancement in ion beam experiments with Bayesian probability theory *Nucl. Instrum. Meth.* **136–138** 1140
- Fischer R, Hanson K M, Dose V and von der Linden W 2000 Background estimation in experimental spectra *Phys. Rev.* **E61** 1152
- Fröhner F H 2000 Evaluation and analysis of nuclear resonance data *JEFF Report 18* (Paris: OECD Publications)
- Gelman A, Carlin J B, Stern H S and Rubin D B 1995 *Bayesian Data Analysis* (London: Chapman and Hall)
- Glimm J and Sharp D H 1999 Prediction and the quantification of uncertainty *Physica D* **133** 152
- Gregory P C and Loredo T J 1992 A new method for the detection of a periodic signal of unknown shape and period *Astron. J.* **398** 146
- Gregory P C and Loredo T J 1996 Bayesian periodic signal detection II—Bayesian periodic signal detection: analysis of ROSAT observations of PSR 0540-693 *Astron. J.* **473** 1059
- Gregory P C 1999 Bayesian periodic signal detection I—analysis of 20 years of radio flux measurements of the x-ray binary LS I +61°303 *Astron. J.* **520** 361
- Gubernatis J E, Jarrell M, Silver R N and Sivia D S 1991 Quantum Monte-Carlo simulations and maximum-entropy: dynamics from imaginary-time data *Phys. Rev. B* **44** 6011
- Hanson K M 1993 Introduction to Bayesian image analysis *Medical Imaging: Image Processing* ed M H Loew *Proc. SPIE* **1898** 716
- Hanson K M 2000 *Tutorial on Markov Chain Monte Carlo* <http://public.lanl.gov/kmh/talks/maxent00b.pdf>
- Hasting W K 1970 Monte Carlo sampling methods using Markov chains and their applications *Biometrika* **57** 97
- Higdon D M and Yamamoto S Y 2001 Estimation of the head sensitivity function in scanning magnetoresistance microscopy *J. Am. Stat. Assoc.* **96** 785
- Hobson M P, Bridle S L and Lahav 2002 Combining cosmological datasets: hyperparameters and Bayesian evidence arXiv:astro-ph/0203259
- Howson C and Urbach P 1993 *Scientific Reasoning—The Bayesian Approach* (Chicago and La Salle: Open Court)
- ISO (International Organization for Standardization) 1993 *Guide to the Expression of Uncertainty in Measurement* (Geneva: ISO)
- Jaynes E T 1957a Information theory and statistical mechanics *Phys. Rev.* **106** 620
- Jaynes E T 1957b Information theory and statistical mechanics II *Phys. Rev.* **108** 171
- Jaynes E T 1968 Prior probabilities *IEEE Trans. Syst. Cybern.* **SSC-4** 227; reprinted in Jaynes 1983
- Jaynes E T 1973 The well-posed problem *Found. Phys.* **3** 477; reprinted in Jaynes 1983
- Jaynes E T 1983 *Papers on Probability, Statistics and Statistical Physics* ed W L Harper and C A Hooker (Dordrecht: Reidel)
- Jaynes E T 1998 <http://omega.albany.edu:8008/JaynesBook.html>
- Jeffreys H 1961 *Theory of Probability* (Oxford: Oxford University Press)
- John M V and Narlikar J V 2002 Comparison of cosmological models using Bayesian theory *Phys. Rev.* **D65** 43506
- Kadane J B and Schum D A 1996 *A Probabilistic Analysis of the Sacco and Vanzetti Evidence* (Chichester: Wiley)
- Kalman R E 1960 A new approach to linear filtering and prediction problems *Trans. ASME J. Casic Eng.* **82** 35
- Kass R E, Carlin B P, Gelman A and Neal R M 1998 Markov chain Monte Carlo in practice: A roundtable discussion *Am. Stat.* **52** 93
- Lad F 1996 *Operational Subjective Statistical Methods—A Mathematical, Philosophical, and Historical Introduction* (Chichester: Wiley)
- Lee P M 1989 *Bayesian statistics—An Introduction* (Chichester: Wiley)
- Lewis A and Bridle S 2002 Cosmological parameters from CMB and other data: a Monte-Carlo approach *Phys. Rev.* **D66** 103511
- von der Linden W 1995 Maximum-entropy data analysis *Appl. Phys.* **A60** 155
- von der Linden W, Dose V and Fischer R 1996b *Proc. 1996 Maximum Entropy Conf. on Spline-Based Adaptive Resolution Image Reconstruction* ed Sears M *et al* (Port Elizabeth: N.M.B. Printers) p 154
- Lindley D V 1986 Discussion to Efron 1986a *Am. Stat.* **40** 6
- Loredo T J 1990 *Maximum Entropy and Bayesian Methods* ed P F Fougère (Dordrecht: Kluwer) p 81

- Loredo T J and Lamb D Q 2002 Bayesian analysis of neutrinos observed from supernova SN 1987A *Phys. Rev.* **D65** 063002
- Malakoff D 1999 Bayes offers a 'New' way to make sense of numbers *Science* **286** 1460
- Maybaeck P S 1979 *Stochastic Models, Estimation and Control* vol 1 (New York: Academic)
- Metropolis H, Rosenbluth A W, Rosenbluth M N, Teller A H and Teller E 1953 Equations of state calculations by fast computing machines *J. Chem. Phys.* **21** 1087
- von Mises R 1957 *Probability, Statistics, and Truth* (St Leonards: Allen and Unwin); reprinted in 1987 by Dover
- Neal R M 1993 *Probabilistic Inference Using Markov Chain Monte Carlo Methods* (Toronto: Technical Report CRG-TR-93-1)
- O'Hagan A 1994 *Kendall's Advanced Theory of Statistics: Vol 2B—Bayesian Inference* (New York: Halsted)
- Pearl J 1988 *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference* (San Mateo: Morgan Kaufmann)
- Press W H 1997 Understanding data better with Bayesian and global statistical methods *Unsolved Problems in Astrophysics* 49–60, ed J N Bahcall and J P Ostriker (Princeton: Princeton University) p 49
- Press S J 2002 *Subjective and Objective Bayesian Statistics: Principles, Models, and Applications* 2nd edn (Chichester: Wiley)
- Robert C P 2001 *The Bayesian Choice* (New York: Springer)
- Saquist S S, Hanson K M and Cunningham G S 1997 Model-based image reconstruction from time-resolved diffusion data *Proc. SPIE* **3034** 369
- Schrödinger E 1947a The Foundation of the theory of probability—I *Proc. R. Irish Acad.* **51A** 51; reprinted in *Collected Papers* vol 1 (Vienna 1984: Austrian Academy of Science) p 463
- Schrödinger E 1947b The Foundation of the theory of probability—II *Proc. R. Irish Acad.* **51A** 141; reprinted in *Collected Papers* vol 1 (Vienna 1984: Austrian Academy of Science) p 479
- Sivia D S 1997 *Data Analysis—a Bayesian Tutorial* (Oxford: Clarendon)
- Skilling J 1992 Quantified maximum entropy *Int. Spectr. Lab.* **2** 4
- Smith A F M 1991 Bayesian numerical analysis *Phys. Trans. R. Soc. London* **337** 369
- Taylor B N and Kuyatt C E 1994 *Guidelines for Evaluating and Expressing Uncertainty of NIST Measurement Results* (Gaithersburg: NIST Technical Note 1297); available on line at <http://physics.nist.gov/>
- Tribus M 1969 *Rational Descriptions, Decisions, and Designs* (Elmsford: Pergamon)
- Welch G and Bishop G 2002 *An Introduction to Kalman Filter* <http://www.cs.unc.edu/~welch/kalman/>
- Zech G 2002 Frequentist and Bayesian confidence limits *Eur. Phys. J. Direct* **C12** 1
- Zellner A 1986 Bayesian solution to a problem posed by Efron *Am. Stat.* **40** 330