

Maximum entropy and Bayesian data analysis: Entropic prior distributions

Ariel Caticha

Physics Department, State University of New York at Albany, Albany, New York 12222, USA

Roland Preuss

Center for Interdisciplinary Plasma Science, Max-Planck-Institut fuer Plasmaphysik, EURATOM Association, Boltzmannstrasse 2, D-85748 Garching bei Muenchen, Germany

(Received 16 July 2003; published 29 October 2004)

The problem of assigning probability distributions which reflect the prior information available about experiments is one of the major stumbling blocks in the use of Bayesian methods of data analysis. In this paper the method of maximum (relative) entropy (ME) is used to translate the information contained in the known form of the likelihood into a prior distribution for Bayesian inference. The argument is inspired and guided by intuition gained from the successful use of ME methods in statistical mechanics. For experiments that cannot be repeated the resulting “entropic prior” is formally identical with the Einstein fluctuation formula. For repeatable experiments, however, the expected value of the entropy of the likelihood turns out to be relevant information that must be included in the analysis. The important case of a Gaussian likelihood is treated in detail.

DOI: 10.1103/PhysRevE.70.046127

PACS number(s): 02.50.Tt, 02.50.Cw, 05.20.-y

I. INTRODUCTION

The inference of physical quantities from data generated either by experiment or by numerical simulation is a ubiquitous and often cumbersome task. Whether the data is corrupted by noise, hampered by finite resolution or tied up in correlations, in principle it should always be possible to improve the analysis by taking into account, in addition to the information contained in the data, whatever other knowledge one might have about the physical quantities to be inferred or about how the data was generated. The way to link this prior information with the new information in the data is found in Bayesian probability theory.

Bayesian methods are increasingly popular in physics [1]. They are essential whenever repeating the experiment many times in order to reduce the measurement uncertainty is either too expensive or time consuming. This is a common situation in astronomy and astrophysics [2], and also in large laboratory experiments as in fusion [3] and in high energy physics [4]. Other typical uses in physics arise in spectrum restoration, in ill-posed inversion problems [5–7] and when separating a signal from an unknown background [8]. Applications include mass spectrometry [9], Rutherford back-scattering [10] and nuclear magnetic resonance [11].

From a general point of view the problem of inductive inference is to update from a prior probability distribution to a posterior distribution when new information becomes available. The challenge is to develop updating methods that are systematic and objective. Two methods have been found which are of very broad applicability: one is based on Bayes’ theorem and the other is based on the maximization of entropy. The choice between these two updating methods is dictated by the nature of the information being processed.

When we want to update our beliefs about the values of certain quantities θ on the basis of the observed values of other quantities y — the data — and of some known relation

between θ and y we must use Bayes’ theorem. The updated or posterior distribution is $p(\theta|y) \propto \pi(\theta)p(y|\theta)$; the relation between y and θ is supplied by a known model $p(y|\theta)$; the previous knowledge about θ is codified both into the “prior” probability $\pi(\theta)$ and also in the “likelihood” distribution $p(y|\theta)$.

The selection of the prior $\pi(\theta)$ is a controversial issue which has generated an enormous literature [12]. The difficulty is that it is not clear how to carry out an objective translation of our previous beliefs about θ into a distribution $\pi(\theta)$. One reasonable attitude is to admit subjectivity and recognize that different individuals may start from the same information and legitimately end with different translations. In simple cases experience and physical intuition have led to a considerable measure of success, but we are often confronted with new complex situations involving perhaps parameter spaces of high dimensionality where we have neither a previous experience nor a reliable intuition.

On the other hand, there are special cases where some degree of objectivity can be attained. One example is provided by those cases where symmetries are known to hold; requirements of invariance can go a long way towards the complete specification of a prior.

A second example is given by Bayes’ theorem itself. Suppose two batches of data, y_1 and y_2 are to be processed. Start from a subjective prior $\pi_1(\theta)$ and use Bayes to process the data y_1 to update to the posterior $p(\theta|y_1) \propto \pi_1(\theta)p(y_1|\theta)$. To process the data y_2 one uses Bayes’ theorem again, $p(\theta|y_1y_2) \propto \pi_2(\theta)p(y_2|\theta y_1)$, where the new prior is given by $\pi_2(\theta) = p(\theta|y_1)$. While there is nothing new in this well-known example of Bayesian updating it is interesting to note that we are confronted with a situation in which information of a very special kind (the data y_1) has been incorporated into a distribution (the prior π_2) in a totally objective way. This raises the question of whether other examples exist where information can be objectively incorporated into a prior dis-

tribution. The purpose of this paper is to offer an affirmative answer by showing explicitly which type of information and which is the appropriate method for processing it.

Other attempts to seek objectivity in assigning priors have sought to characterize that elusive state of knowledge which presumably reflects complete ignorance. Although there are convincing arguments against the existence of such noninformative priors [13], the search has had the merit of suggesting connections with the notion of entropy [14] including two proposals for “entropic priors” [15,16]. This brings us to the second method of processing information.

Bayes’ theorem follows from the product rule for joint probabilities, $p(y, \theta) = \pi(\theta)p(y|\theta)$, and therefore its applicability is restricted to situations where assertions concerning the joint values of the data y and the parameters θ are meaningful. There are, however, situations where the available information is of a different nature and involves assertions about the probabilities themselves. Such information, which includes but is not limited to assertions about expected values, cannot be processed using Bayes’ theorem.

The method of maximum entropy (ME) is designed for updating from a prior probability distribution to a posterior distribution when the information to be processed takes the form of a constraint on the family of acceptable posterior distributions [17]. The early and less satisfactory justification of the ME method followed from interpreting entropy, through the Shannon axioms, as a measure of the amount of uncertainty in a probability distribution [18,19]. Objections to this approach are that the Shannon axioms refer to probabilities of discrete variables, the entropy of continuous variables is not defined, and that the use of one particular entropy as the unique measure of uncertainty remained questionable. Other so-called entropies could and, indeed, were introduced. Ultimately, the real problem is that Shannon was not concerned with inductive inference. He was not trying to update probability distributions but was instead analyzing the capacity of communication channels. Shannon’s entropy makes no reference to prior distributions.

Considerations such as these motivated several attempts to justify the ME method directly as a method of inductive inference without invoking questionable measures of uncertainty [20–22]. In these approaches the concept of relative entropy is then introduced as a tool for consistent reasoning which, in the special case of uniform priors, reduces to the usual entropy. There is no need for an interpretation in terms of heat, disorder, or uncertainty, or even in terms of an amount of information. Perhaps this is the explanation of why the search for the meaning of entropy has turned out to be so elusive: strictly, *entropy needs no interpretation*. In Sec. II, as background for the rest of the paper, we present a brief outline of one such no-interpretation approach which is more extensively developed elsewhere [23].

In this paper we use entropic arguments to translate prior information into a prior distribution. Rather than seeking a totally noninformative prior, we make use of information that we do in fact have. Remarkably, it turns out that the very condition that allows us to contemplate using Bayes’ theorem — namely, the knowledge of a likelihood function, $p(y|\theta)$ — already constitutes valuable prior information. In this sense one can assert that the search for completely noninformative

priors is misplaced: if we do not know the likelihood, then prior distributions are not needed anyway. The prior thus obtained is an “entropic prior.” The name and the first proposal of a prior of this kind is due to Skilling [15] for the case of discrete distributions. The generalization to the continuous case and further elaborations by Rodríguez [16,25] constitute a second proposal.

It is essential for the successful use of any prior, and of entropic priors in particular, to be aware of what information they contain and, crucially, what information they do not contain. No prior can be expected to succeed unless all the information relevant to the problem at hand has been taken into account. It is quite likely that most practical problems that were encountered with entropic priors in the past can be traced to a failure to identify and incorporate all the relevant information.

The information that has, in this paper, been translated into the entropic prior is that contained in the likelihood. The *bare* entropic priors discussed here apply to a situation where all we know about the quantities θ is that they appear as parameters in the likelihood $p(y|\theta)$, and *nothing else*. Generalizations are, of course, possible. Sometimes we are aware of additional relevant information beyond what is contained in the likelihood and it can easily be incorporated into a modified entropic prior. Other times we might be guilty of overlooking additional information we already have. Indeed, we would not be willing to spend valuable effort in the determination of a parameter θ unless we suspected that knowledge of θ has important implications elsewhere. Typically we know something about the physical significance and the physical meaning of θ . It is clear that in these cases we know considerably more than just that θ is a parameter appearing in the likelihood. We might even conceive of several different experiments, $e=1,2,\dots$, each yielding different sets of data y_e related to θ by different likelihood functions $p_e(y_e|\theta)$. It is sometimes objected that one’s prior knowledge about θ should not depend on which experiment one decides to use to measure it, but this objection is misplaced: the mere fact that θ is measurable through one or another experiment is additional information which, if relevant, should be taken into account.

Another family of problems that can be tackled as a rather straightforward extension of the ideas described here involve choosing which likelihood distribution from among several competing candidates is responsible for generating the data. Indeed, it is clear that any systematic approach to model selection requires as a prerequisite the capability to process in an objective way the information implicit in each of those likelihoods. Except for some brief remarks in the final section, all these further developments, valuable as they might be, will be addressed elsewhere.

Our contribution includes a derivation of an entropic prior (Sec. III) following the same principles of ME inference that have been successful in statistical mechanics. In fact, our whole approach is guided by intuition gained from applications of ME to statistical mechanics. Preliminary steps along this direction were taken in Ref. [26] where a problem with the important case of experiments that can be indefinitely repeated had already been identified but not fully resolved. This problem, re-examined in Sec. IV, is interpreted as a

symptom that important relevant information has been overlooked. The complete resolution, which hinges on identifying and incorporating this additional information, is given in Secs. V and VI. The actual way in which ME is used in the derivation, in analogy to standard applications in statistical mechanics, turns out to be important because it clarifies what it is that has been derived and how to use it: ours is, in effect, a third proposal for an entropic prior. In Sec. VII we discuss in detail the important example of a Gaussian likelihood and finally, in Sec. VIII, we summarize and comment on the differences among the three versions of entropic prior and on possible further developments.

II. THE LOGIC BEHIND THE ME METHOD

The goal is to update beliefs about $y \in Y$ which are codified in the prior probability distribution $m(y)$ to a posterior distribution $p(y)$ when new information in the form of a constraint becomes available. (The constraints can, but need not, be linear.) The selection is carried out by ranking the probability distributions according to increasing *preference*. One feature we impose on the ranking scheme is transitivity: if distribution p_1 is preferred over distribution p_2 , and p_2 is preferred over p_3 , then p_1 is preferred over p_3 . Such transitive rankings are implemented by assigning to each $p(y)$ a real number $S[p]$ called the entropy of p in such a way that if p_1 is preferred over p_2 , then $S[p_1] > S[p_2]$. The selected p will be that which maximizes $S[p]$. Thus the method involves entropies which are real numbers and entropies that should be maximized. These are features imposed by design; they are dictated by the function that the ME method is supposed to perform.

Next we determine the functional form of $S[p]$. This is the rule that defines the ranking scheme. *The purpose of the rule is to do induction.* We want to extrapolate, to generalize from those special cases where we know what the preferred distribution should be to the much larger number of cases where we do not. Thus, in order to be an inductive rule $S[p]$ must have wide applicability; we will assume that *the same rule applies to all cases.* There is no justification for this universality except for the usual pragmatic justification of induction: we must be inclined to generalize lest we become paralyzed into not generalizing at all. But then, we should remain cautious and keep in mind that in many instances induction just fails.

The argument goes as follows [22]. If a general theory exists, then it must apply to special cases. Furthermore, if in a certain special case the preferred distribution is known, then this knowledge can be used to constrain the form of $S[p]$. Finally, if enough special cases are known, then $S[p]$ will be completely determined. The known special cases are called the “axioms” of ME. As we will see below the axioms reflect the conviction that one should not change one’s mind frivolously, that whatever was learned in the past is important. The chosen posterior distribution should coincide with the prior as closely as possible and one should only update those aspects of one’s beliefs for which corrective new evidence has been supplied. The three axioms are listed below.

Axiom 1: Locality. *Local information has local effects.*

When the new information does not refer to a domain $D \subset Y$, then the conditional probabilities $p(y|D)$ will not be revised. The power of this axiom stems from the arbitrariness in the choice of D . The consequence of the axiom is that nonoverlapping domains of y contribute additively to the entropy: $S[p] = \int dy F[p(y)]$ where F is some unknown function.

Axiom 2: Coordinate invariance. *The ranking should not depend on the system of coordinates.* The coordinates that label the points y are arbitrary; they carry no information. The consequence of this axiom is that $S[p] = \int dy p(y)f[p(y)/m(y)]$ involves coordinate invariants such as $dy p(y)$ and $p(y)/m(y)$, where the density $m(y)$ and the function f are, at this point, unknown.

Next we make a second use of the locality axiom to enforce objectivity. We allow domain D to extend over the whole space Y and assert that *when there is no new information there is no reason to change one’s mind.* When there are no constraints the selected posterior distribution should coincide with the prior distribution. This eliminates the arbitrariness in the density $m(y)$: up to normalization $m(y)$ is the prior distribution.

Axiom 3: Consistency for independent subsystems. *When a system is composed of subsystems that are believed to be independent it should not matter whether the inference procedure treats them separately or jointly.* If we originally believe that two systems are independent and the new evidence is silent on the matter of correlations, then there is no reason to change our mind. Specifically, if $y = (y_1, y_2) \in Y = Y_1 \times Y_2$, and the subsystem priors $m_1(y_1)$ and $m_2(y_2)$ are, respectively, upgraded to $p_1(y_1)$ and $p_2(y_2)$, then the prior for the whole system $m_1(y_1)m_2(y_2)$ should be upgraded to $p_1(y_1)p_2(y_2)$. This axiom restricts the function f to be a logarithm. (The fact that the logarithm applies also when the subsystems are not independent follows from our inductive hypothesis that the ranking scheme has universal applicability.)

The overall consequence of these axioms [27] is that probability distributions $p(y)$ should be ranked relative to the prior $m(y)$ according to their (relative) entropy [17],

$$S[p, m] = - \int dy p(y) \log \frac{p(y)}{m(y)}. \tag{1}$$

The derivation has singled out $S[p, m]$ as *the unique entropy to be used in inductive inference.* Other expressions, such as $S[m, p]$, or $S[p, m] + S[m, p]$, or even expressions that do not involve the logarithm, may be useful for other purposes, but they do not constitute an induction: they are not a generalization from the simple cases described in the axioms.

We end this section with two comments on the prior density $m(y)$. First, $S[p, m]$ may be infinitely negative when $m(y)$ vanishes within some region D . In other words, the ME method confers an overwhelming preference on those distributions $p(y)$ that vanish whenever $m(y)$ does. Is this a problem? Not really. A similar “problem” also arises in the context of Bayes’ theorem. A vanishing prior represents a tremendously serious prejudice because no amount of data to the contrary would allow us to revise it. The solution in both cases is to recognize that unless we are absolutely certain that y could not possibly lie within D then we should not

have assigned $m(y)=0$ in the first place. Assigning a very low but nonzero prior represents a safer and less prejudiced representation of one's beliefs and/or doubts both in the context of Bayesian and of ME inference.

Second, choosing the prior density $m(y)$ can be tricky. When there is no information leading us to prefer one microstate of a physical system over another we might as well assign equal prior probability to each state. Thus it is reasonable to identify $m(y)$ with the density of states and the invariant $m(y)dy$ is the number of microstates in dy . This is the basis for statistical mechanics. Other examples of relevance to physics arise when there is no reason to prefer one region of the space Y over another. Then we should assign the same prior probability to regions of the same "volume," and we can choose $\int_R dy m(y)$ to be the volume of a region R in the space Y . Notice that because of the presence of the prior $m(y)$ not all subjectivity has been eliminated and Laplace's principle of insufficient reason still plays an important role, albeit in a somewhat modified form. Just as with Bayes' theorem, what is objective here is the manner in which information is processed, not the initial probability assignments.

III. ENTROPIC PRIORS: THE BASIC IDEA

In this section we follow [26] closely. We use the ME method to derive a prior $\pi(\theta)$ for use in Bayes' theorem,

$$p(\theta|y) \propto p(y, \theta) = \pi(\theta)p(y|\theta). \quad (2)$$

The selection of a preferred distribution using the ME method demands that the space in which the search will be conducted be specified. Being a consequence of the product rule for joint probabilities, Bayes' theorem requires that assertions such as y and θ be meaningful and that the probability of y and θ be well defined. Therefore we must focus our attention on $p(y, \theta)$ rather than $\pi(\theta)$; the relevant universe of discourse is neither Θ , the space of all θ s, nor the data space Y , but the product $\Theta \times Y$. This point, first made by Rodríguez [24], is central to the argument. Our derivation and the final result, however, differ from his [24,25] in several respects.

To rank distributions in the space $\Theta \times Y$ we must decide on a prior $m(y, \theta)$. At this starting point absolutely nothing is known about the variables θ , in particular, they have no physical meaning, and no relation between y and θ is known. The θ s are totally arbitrary. Therefore the prior must be a product $m(y)\mu(\theta)$ of the separate priors in the spaces Y and Θ . Indeed, the distribution that maximizes the relative entropy,

$$\sigma[p] = - \int dy d\theta p(y, \theta) \log \frac{p(y, \theta)}{m(y)\mu(\theta)}, \quad (3)$$

when no constraints are imposed is $p(y, \theta) \propto m(y)\mu(\theta)$; it is such that data about y tells us absolutely nothing about θ .

In what follows we assume that $m(y)$ is known. We consider this an important part of understanding what data it is that has been collected. In Sec. VII we will suggest a reasonable $m(y)$ for the special case of a Gaussian likelihood. The prior $\mu(\theta)$ remains unspecified.

Next we incorporate the crucial piece of information from which the parameters θ derive their physical meaning and which establishes the relation between θ and y : the likelihood function $p(y|\theta)$ is known. This has two consequences: First, the joint distribution $p(y, \theta)$ is constrained to be of the form $\pi(\theta)p(y|\theta)$. Notice that this constraint is not in the form that is most usual for applications of the ME method: it is not an expectation value. Note also that the only information we are using about the quantities θ is that they appear as parameters in the likelihood $p(y|\theta)$, *nothing else*. In many situations of experimental interest there exists additional relevant information beyond what is contained in the likelihood; such information should be included as additional constraints in the maximization of the relative entropy σ .

Second, now that a bare minimum is known about θ , namely that each θ represents a probability distribution, there is a natural but still subjective choice for $\mu(\theta)$. As discussed in Ref. [28], except for an overall multiplicative constant, there is a unique Riemannian metric that adequately reflects the fact that the points in a space of probability distributions are not structureless, but happen to be probability distributions; this is the Fisher-Rao metric. Within the finite-dimensional subspace defined by the constraint — the known $p(y|\theta)$ — the natural metric on Θ is $d\ell^2 = g_{ij}d\theta^i d\theta^j$, where the unique g_{ij} induced by the family of distributions $p(y|\theta)$ is

$$g_{ij} = \int dy p(y|\theta) \frac{\partial \log p(y|\theta)}{\partial \theta^i} \frac{\partial \log p(y|\theta)}{\partial \theta^j}. \quad (4)$$

Accordingly we choose $\mu(\theta) = g^{1/2}(\theta)$, where $g(\theta)$ is the determinant of g_{ij} . Having identified the prior measure and the constraints, we allow the ME method to take over.

The preferred distribution $p(y, \theta)$ is chosen by varying $\pi(\theta)$ to maximize

$$\begin{aligned} \sigma[\pi] &= - \int dy d\theta \pi(\theta) p(y|\theta) \log \frac{\pi(\theta)p(y|\theta)}{g^{1/2}(\theta)m(y)} \\ &= - \int d\theta \pi(\theta) \log \frac{\pi(\theta)}{g^{1/2}(\theta)} + \int d\theta \pi(\theta) S(\theta), \end{aligned} \quad (5)$$

where $S(\theta)$ is the entropy of the likelihood,

$$S(\theta) = - \int dy p(y|\theta) \log \frac{p(y|\theta)}{m(y)}. \quad (6)$$

Writing the Langrange multiplier that enforces $\int d\theta \pi(\theta) = 1$ as $1 - \log \zeta$, and assuming $p(y|\theta)$ is normalized yields

$$0 = \int d\theta \left(- \log \frac{\pi(\theta)}{g^{1/2}(\theta)} + S(\theta) - \log \zeta \right) \delta\pi(\theta). \quad (7)$$

Therefore, the probability that the value of θ should lie within the small volume $g^{1/2}(\theta)d\theta$ is

$$\pi(\theta)d\theta = \frac{1}{\zeta} e^{S(\theta)} g^{1/2}(\theta) d\theta \quad \text{with} \quad \zeta = \int d\theta g^{1/2}(\theta) e^{S(\theta)}. \quad (8)$$

This entropic prior is our first main result. It tells us that the preferred value of θ is that which maximizes the entropy

$S(\theta)$ because this maximizes the scalar probability density exp $S(\theta)$. It also tells us the degree to which values of θ away from the maximum are ruled out; in many cases the preference for the ME distribution can be overwhelming. Note also that the density exp $S(\theta)$ is a scalar function and the presence of the Jacobian factor $g^{1/2}(\theta)$ makes Eq. (8) manifestly invariant under changes of the coordinates θ in the space Θ .

We can claim a partial success. The ingredients that have been used are precisely those that led us to consider using Bayes' theorem in the first place. The information contained in the model — by which we mean that the data space Y , its measure $m(y)$, and the conditional distribution $p(y|\theta)$ — has been translated into a prior $\pi(\theta)$. The success is partial because it has been achieved for the special case of a fixed data space Y of experiments that cannot be repeated. A more complete treatment requires that we address the important case of experiments that can be repeated indefinitely.

IV. REPEATABLE EXPERIMENTS

Experiments need not be repeatable but sometimes they are. Let us assume that successive repetitions are possible and that they happen to be independent. Suppose, to be specific, that the experiment is performed twice so that the space of data $Y \times Y = Y^2$ consists of the possible outcomes y_1 and y_2 . Suppose further that θ is not a “random” variable; the value of θ is fixed but unknown. Then the joint distribution in the space $\Theta \times Y^2$ is

$$p(y_1, y_2, \theta) = \pi^{(2)}(\theta) p(y_1, y_2 | \theta) = \pi^{(2)}(\theta) p(y_1 | \theta) p(y_2 | \theta), \quad (9)$$

and the appropriate σ entropy is

$$\begin{aligned} \sigma^{(2)}[\pi] = & \\ & - \int dy_1 dy_2 d\theta p(y_1, y_2, \theta) \log \frac{p(y_1, y_2, \theta)}{[g^{(2)}(\theta)]^{1/2} m(y_1) m(y_2)}, \end{aligned} \quad (10)$$

where $g^{(2)}(\theta)$ is the determinant of the Fisher-Rao metric for $p(y_1, y_2 | \theta)$. From Eq. (4) it follows that $g_{ij}^{(2)} = 2g_{ij}$ so that $g^{(2)}(\theta) = 2^d g(\theta)$, d being the dimension of θ . Maximizing $\sigma^{(2)}[\pi]$ subject to $\int d\theta \pi^{(2)}(\theta) = 1$ we get

$$\pi^{(2)}(\theta) = \frac{1}{Z^{(2)}} g^{1/2}(\theta) e^{S^{(2)}(\theta)} = \frac{1}{Z^{(2)}} g^{1/2}(\theta) e^{2S(\theta)}, \quad (11)$$

where $S^{(2)}(\theta) = 2S(\theta)$ is the entropy of $p(y_1, y_2 | \theta)$, $S(\theta) \stackrel{\text{def}}{=} S^{(1)}(\theta)$ and $Z^{(2)}$ is the appropriate normalization constant. The generalization to N repetitions of the experiment, with data space Y^N , is immediate,

$$\pi^{(N)}(\theta) = \frac{1}{Z^{(N)}} g^{1/2}(\theta) e^{S^{(N)}(\theta)} = \frac{1}{Z^{(N)}} g^{1/2}(\theta) e^{NS(\theta)}. \quad (12)$$

This is clearly wrong: the dependence of $\pi^{(N)}$ on the amount N of data would lead us to a perpetual revision of the prior as more data is collected. The absurdity of this situation be-

comes manifest when we consider the case of large N . Then the exponential preference for the value of θ that maximizes $S(\theta)$ becomes so pronounced that no amount of data to the contrary can successfully overcome its effect. The data becomes irrelevant, and the more data we have, the more irrelevant it becomes.

Repeatable experiments present us with a problem. One possible attitude is to blame the ME method: it gives nonsense and cannot be trusted. As with all inductive methods this is, of course, a logical possibility. A second, more constructive approach, is to always be prepared to question the results of ME calculations on the basis that there is no guarantee that all the information relevant to the situation at hand has been taken into account. The problem is not a failure of the ME method but a failure to include all the relevant information.

That this is indeed the case can be seen as follows: When we say an experiment can be repeated twice, $N=2$, we actually know more than just $p(y_1, y_2 | \theta) = p(y_1 | \theta) p(y_2 | \theta)$. We also know that forgetting or discarding the value of say y_2 , yields an experiment that is totally indistinguishable from the single, $N=1$, experiment. This *additional* information is quantitatively expressed by $\int dy_2 p(y_1, y_2, \theta) = p(y_1, \theta)$, or equivalently

$$\int dy_2 \pi^{(2)}(\theta) p(y_1 | \theta) p(y_2 | \theta) = \pi^{(1)}(\theta) p(y_1 | \theta), \quad (13)$$

which leads to $\pi^{(2)}(\theta) = \pi^{(1)}(\theta)$. In the general case we get the manifestly reasonable result

$$\pi^{(N)}(\theta) = \pi^{(N-1)}(\theta) = \dots = \pi^{(1)}(\theta). \quad (14)$$

The challenge then is to identify a constraint that codifies this information within each space $\Theta \times Y^N$.

V. MORE INFORMATION: THE LAGRANGE MULTIPLIER α

The problem with the prior $\pi^{(N)}(\theta)$ in Eq. (12) is that it expresses an overwhelming preference for the value θ_{\max} of θ that maximizes the entropy $S(\theta)$. Indeed, as $N \rightarrow \infty$ we have $\pi^{(N)}(\theta) \rightarrow \delta(\theta - \theta_{\max})$ leading to

$$\langle S \rangle = \int d\theta \pi^{(N)}(\theta) S(\theta) \xrightarrow{N \rightarrow \infty} S(\theta_{\max}). \quad (15)$$

This suggests that a better prior would be obtained by maximizing the entropy $\sigma^{(N)}$ of distributions on the space $\Theta \times Y^N$ subject to an additional constraint on the numerical value \bar{S} of the expected entropy $\langle S \rangle$. It is not that we happen to know the numerical value \bar{S} of $\langle S \rangle$; in fact we do not. It is rather that we recognize that information about \bar{S} is relevant in the sense that if \bar{S} were known the problem above would not arise. Naturally, additional effort will be required to obtain the needed value of \bar{S} .

The logic of the previous paragraph may sound unfamiliar and further comments may be helpful. When justifying the use of the ME method to obtain, say, the canonical

Boltzmann-Gibbs distribution ($P_q \propto e^{-\beta E_q}$) it has been common to say something like “we seek the minimally biased (i.e., maximum entropy) distribution that codifies the information we do possess (the expected energy) and nothing else.” Many authors find this justification objectionable. Indeed, they might argue, for example, that the spectrum of black body radiation is what it is independently of whatever information happens to be available to us. We prefer to phrase our objection differently: in most realistic situations the expected value of the energy is not a quantity we happen to know. Nevertheless, it is still true that maximizing entropy subject to a constraint on this (unknown) expected energy leads to correct predictions. Therefore, the justification behind imposing a constraint on the expected energy cannot be that this is a quantity that happens to be known — because it is not — but rather that the expected energy is the quantity that *should* be known. Even if unknown, we recognize it as the crucial relevant information without which no successful predictions can be made. Therefore we proceed as if this crucial information were available and produce a formalism that contains the temperature as a free parameter that will later have to be obtained from the experiment itself. In other words, the temperature (or expected energy) is one additional parameter to be inferred from the data.

The entropy on the space $\Theta \times Y^N$ is

$$\begin{aligned} \sigma^{(N)}[\pi] = & - \int d\theta dy^{(N)} \pi(\theta) p(y^{(N)}|\theta) \log \frac{\pi(\theta) p(y^{(N)}|\theta)}{g^{1/2}(\theta) m(y^{(N)})} = \\ & - \int d\theta \pi(\theta) \log \frac{\pi(\theta)}{g^{1/2}(\theta)} + N \int d\theta \pi(\theta) S(\theta), \end{aligned} \quad (16)$$

where $S(\theta)$ given by Eq. (6). [A constant factor of $N^{d/2}$ associated to the Fisher-Rao measure $g^{(N)}(\theta)$ has been omitted. It would eventually be absorbed into the normalization of $\pi(\theta)$.] To obtain the prior $\pi(\theta)$ we maximize $\sigma^{(N)}$ subject to constraints on $\langle S \rangle$ and that π be normalized,

$$\begin{aligned} \delta \left[\sigma^{(N)} + (1 - \log \zeta) \left(\int d\theta \pi(\theta) - 1 \right) \right. \\ \left. + \lambda_N \left(\int d\theta \pi(\theta) S(\theta) - \bar{S} \right) \right] = 0. \end{aligned} \quad (17)$$

This gives,

$$\int d\theta \left(- \log \frac{\pi(\theta)}{g^{1/2}(\theta)} + (N + \lambda_N) S(\theta) - \log \zeta \right) \delta\pi(\theta) = 0. \quad (18)$$

Therefore,

$$\pi(\theta) = \frac{1}{\zeta} g^{1/2}(\theta) \exp[(N + \lambda_N) S(\theta)]. \quad (19)$$

The undesired dependence on N is eliminated if in each space $\Theta \times Y^N$ the Lagrange multipliers λ_N are chosen so that $N + \lambda_N = \alpha$ is a constant independent of N . The resulting entropic prior,

$$\pi(\theta|\alpha) = \frac{1}{\zeta(\alpha)} g^{1/2}(\theta) e^{\alpha S(\theta)}, \quad (20)$$

satisfies Eq. (14). This is our second main result. The prior $\pi(\theta|\alpha)$ codifies information contained in the likelihood function, plus information about the expected value of the entropy of the likelihood implicit in the hyperparameter α ,

$$\bar{S}(\alpha) = \frac{d}{d\alpha} \log \zeta(\alpha), \quad (21)$$

with $\zeta(\alpha)$ is given by

$$\zeta(\alpha) = \int d\theta g^{1/2}(\theta) e^{\alpha S(\theta)}. \quad (22)$$

The next and final step is to figure out which α applies to the particular experimental situation under consideration. The natural way to proceed is to invoke Bayes' theorem

$$p(\alpha, \theta|y^N) = \pi(\alpha, \theta) \frac{p(y^N|\theta)}{p(y^N)} = \pi(\alpha) \pi(\theta|\alpha) \frac{p(y^N|\theta)}{p(y^N)}. \quad (23)$$

The choice of a prior $\pi(\alpha)$ for α itself is addressed in the next section. If we were truly interested in the actual α , we could marginalize over θ to obtain

$$p(\alpha|y^N) = \int d\theta p(\alpha, \theta|y^N) = \frac{\pi(\alpha)}{p(y^N)} \int d\theta \pi(\theta|\alpha) p(y^N|\theta). \quad (24)$$

But our interest in the value of α is only indirect; α is a necessary but annoying technical complication that stands in our way to the real goal of inferring θ . Marginalizing over α , we get

$$p(\theta|y^N) = \int d\alpha p(\alpha, \theta|y^N) = \bar{\pi}(\theta) \frac{p(y^N|\theta)}{p(y^N)}, \quad (25)$$

where

$$\bar{\pi}(\theta) = \int d\alpha \pi(\alpha) \pi(\theta|\alpha). \quad (26)$$

This is the answer we sought: the effective prior for θ , the averaged $\bar{\pi}(\theta)$, is independent of the actual data y^N , as it should. The last step is the assignment of $\pi(\alpha)$.

VI. AN ENTROPIC PRIOR FOR α

To remain consistent with the spirit of this paper, namely using ME to obtain priors, the prior for α must itself be an entropic prior. The motivation behind discussing entropic priors is that we wish to consider information included in the likelihood function. Since $p(y|\theta)$ refers to θ but makes no reference to any hyperparameters it is quite clear that α should not be treated like the other θ s. The relation between α and the data y is indirect: α is related to θ , and θ is related to y . Once θ is given, the data y becomes irrelevant, it contains no further information about α . The whole significance

of α is derived purely from its appearance in $\pi(\theta|\alpha)$, Eq. (20). Therefore, the relevant universe of discourse is $A \times \Theta$ with $\alpha \in A$. We focus our attention on the joint distribution

$$\pi(\alpha, \theta) = \pi(\alpha)\pi(\theta|\alpha) \quad (27)$$

and we obtain $\pi(\alpha)$ by maximizing the entropy

$$\Sigma[\pi] = - \int d\alpha d\theta \pi(\alpha, \theta) \log \frac{\pi(\alpha, \theta)}{\gamma^{1/2}(\alpha)g^{1/2}(\theta)}, \quad (28)$$

where $\gamma^{1/2}(\alpha)$ is determined below. Since no reference is made to repeatable experiments in Y^N there is no need for any further constraints except for normalization.

The Fisher-Rao measure $\gamma^{1/2}(\alpha)$ in Eq. (28) is

$$\gamma(\alpha) = \int d\theta \pi(\theta|\alpha) \left[\frac{d}{d\alpha} \log \pi(\theta|\alpha) \right]^2. \quad (29)$$

Using Eqs. (20)–(22) we get

$$\gamma(\alpha) = \int d\theta \pi(\theta|\alpha) \left[S(\theta) - \frac{d \log \zeta(\alpha)}{d\alpha} \right]^2 = (\Delta S)^2, \quad (30)$$

but

$$\begin{aligned} \frac{d\bar{S}(\alpha)}{d\alpha} &= \frac{d}{d\alpha} \frac{1}{\zeta(\alpha)} \frac{d\zeta(\alpha)}{d\alpha} = \frac{1}{\zeta(\alpha)} \frac{d^2\zeta(\alpha)}{d\alpha^2} - \left[\frac{1}{\zeta(\alpha)} \frac{d\zeta(\alpha)}{d\alpha} \right]^2 \\ &= (\Delta S)^2. \end{aligned} \quad (31)$$

Therefore,

$$\gamma(\alpha) = \frac{d^2 \log \zeta(\alpha)}{d\alpha^2}. \quad (32)$$

The interpretation is straightforward: the distance between $\pi(\theta|\alpha)$ and $\pi(\theta|\alpha+d\alpha)$ is given by

$$\gamma^{1/2}(\alpha)d\alpha = \Delta S(\alpha)d\alpha, \quad (33)$$

or, in words, the local entropy uncertainty ΔS is the distance per unit change in α .

To maximize Σ rewrite it as

$$\Sigma[\pi] = - \int d\alpha \pi(\alpha) \log \frac{\pi(\alpha)}{\gamma^{1/2}} + \int d\alpha \pi(\alpha) s(\alpha), \quad (34)$$

where $s(\alpha)$ is given by

$$s(\alpha) = - \int d\theta \pi(\theta|\alpha) \log \frac{\pi(\theta|\alpha)}{g^{1/2}(\theta)} = \log \zeta(\alpha) - \alpha \frac{d \log \zeta(\alpha)}{d\alpha}. \quad (35)$$

Then, varying with respect to $\pi(\alpha)$ gives

$$\pi(\alpha) = \frac{1}{z} \gamma^{1/2}(\alpha) e^{s(\alpha)} \quad \text{with } z = \int d\alpha \gamma^{1/2}(\alpha) e^{s(\alpha)}. \quad (36)$$

This completes our derivation.

To summarize: the actual entropic prior for θ in Bayes' theorem, Eq. (25), is $\bar{\pi}(\theta)$ obtained by averaging $\pi(\alpha, \theta)$ over α , Eq. (26), and using Eqs. (36), (20), and (6). The prior $\bar{\pi}(\theta)$ codifies two pieces of information. The first piece is the information contained in the functional form of the likelihood. The second is the expected entropy of the likelihood, which, even when unavailable, turns out to be highly relevant information for repeatable experiments. The actual derivation consists of maximizing the entropy on the space $\Theta \times Y^N$, Eq. (16), subject to constraints corresponding to each of those two pieces of information.

We argued above that the hyperparameter α should not be treated in the same way as the other parameters θ because the likelihood $\pi(y|\theta)$ refers only to θ s and not to α . Nonetheless, it may still be worthwhile to discuss briefly what would happen if α were treated as one of the θ s. In this case, the entropic prior $\pi(\alpha)$ would be determined by focusing our attention on the joint distribution

$$p(\alpha, \theta, y^N) = \pi(\alpha)\pi(\theta|\alpha)p(y^N|\theta), \quad (37)$$

where the last two factors on the right-hand side are assumed known. The assumed universe of discourse would be $A \times \Theta \times Y^N$. A straightforward application of the ME method would, as before, run into trouble with an unwanted N dependence which would require the introduction of a new constraint on the appropriate expected entropy. Thus, the entropic prior for α would involve a second hyperparameter α_2 . The unknown α_2 would itself require its own entropic prior, involving yet a third hyperparameter α_3 , and so on. There would be an endless chain of hyperparameters [16]. In any practical calculation, the chain would have to be truncated. Whether the predictions about θ depend on where and how the truncation is carried out remains to be studied. But, fortunately, this is not necessary: α is not like the other θ s.

VII. EXAMPLE: A GAUSSIAN MODEL

Consider data $y^N = \{y_1, \dots, y_N\}$ that are scattered around an unknown value μ ,

$$y = \mu + \nu \quad (38)$$

with $\langle \nu \rangle = 0$ and $\langle \nu^2 \rangle = \sigma^2$. The goal is to estimate the parameters $\theta = (\theta^1, \theta^2) = (\mu, \sigma)$ on the basis of the data y^N and the information implicit in the model: the data space Y , the measure $m(y)$ (discussed below), and the Gaussian likelihood,

$$p(y|\mu, \sigma) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left[- \frac{(y - \mu)^2}{2\sigma^2} \right]. \quad (39)$$

In Sec. III we asserted that knowing the measure $m(y)$ is part of knowing what data has been collected. Therefore, nothing can be said about $m(y)$ without further specification of the experimental situation. For example, in discrete problems such as tossing a die the symmetry leads to the measure $m(i) = 1$; in the statistical mechanics of discrete energy levels we choose $m(i)$ equal to the degeneracy of the levels; if the data are points in a continuous space it may be appropriate to choose $m(y)dy$ to be a volume element, and if the space is curved then $m(y)$ is nontrivial.

It turns out, however, that in many physical situations where the data happen to be distributed according to a Gaussian the underlying space Y is sufficiently symmetric, e.g., invariant under translations, that we can assume $m(y) = m = \text{constant}$. This is physically reasonable. Gaussian distributions arise when the measured value of y is the sum of a large number of “microscopic” contributions and the details of how the individual contributions are themselves distributed are washed out in the “macroscopic” sum. The macroscopically relevant features are just those that distinguish one Gaussian from another, namely, the mean μ and the variance σ^2 . This is the physical basis behind the Central Limit Theorem. But if microscopic details are irrelevant it should be possible to understand the situation from a purely macroscopic point of view: it should be possible to obtain the Gaussian distribution as the preferred one among all those with the given μ and σ^2 , and this is, indeed, the case: setting $m(y) = \text{constant}$ in $S[p, m]$, Eq. (1), and maximizing subject to constraints on the mean and variance yields Eq. (39).

From Eqs. (6) and (39) the entropy of the likelihood is

$$S(\mu, \sigma) = \log \left[\frac{\sigma}{\sigma_0} \right] \text{ where } \sigma_0 \stackrel{\text{def}}{=} \left(\frac{e}{2\pi} \right)^{1/2} \frac{1}{m}, \quad (40)$$

and the corresponding Fisher-Rao measure, from Eq. (4) is

$$g(\mu, \sigma) = \det \begin{vmatrix} 1/\sigma^2 & 0 \\ 0 & 2/\sigma^2 \end{vmatrix} = \frac{2}{\sigma^4}. \quad (41)$$

Note that both $S(\mu, \sigma)$ and $g(\mu, \sigma)$ are independent of μ . This means that if we were concerned with the simpler problem of estimating μ in a situation where σ happens to be known, then the entropic prior, in any of the versions, Eqs. (8) and (20), or (26), is a constant independent of μ . In other words, when σ is known, the Bayesian estimate of μ using entropic priors coincides with the maximum likelihood estimate, i.e., by the popular procedure of minimizing

$$\chi^2 = \frac{1}{\sigma^2} \sum_{i=1}^N (y_i - \mu)^2. \quad (42)$$

Returning to the more interesting case of unknown σ , the α -dependent entropic prior, Eq. (20) is

$$\pi(\mu, \sigma | \alpha) = \frac{2^{1/2} \sigma^{\alpha-2}}{\zeta(\alpha) \sigma_0^\alpha}. \quad (43)$$

$\pi(\mu, \sigma | \alpha)$ is improper in both μ and σ ; normalization requires the introduction of high and low cutoffs for both μ and σ . The fact that without cutoffs the model is not well defined is an indication that more relevant information is being requested: the cutoffs constitute relevant information that must be taken into account. (The logic parallels that which led to the introduction of α in Sec. V.) The case of unknown cutoff values is important and we intend to explore it in detail in future work. The basic idea is that specifying cutoffs is an integral part of defining the model, and therefore the choice of cutoffs can be tackled as a problem of

model selection. In the remainder of this section, however, we will assume that the information about cutoffs is already available.

It is convenient to write the range of μ as $\Delta\mu = \mu_H - \mu_L$ and to define the σ cutoffs in terms of dimensionless quantities ε_L and ε_H ; σ extends from $\sigma_L = \sigma_0 \varepsilon_L$ to $\sigma_H = \sigma_0 / \varepsilon_H$. Then $\zeta(\alpha)$ and $\pi(\mu, \sigma | \alpha)$ are given by

$$\zeta(\alpha) = \frac{2^{1/2} \Delta\mu \varepsilon_H^{1-\alpha} - \varepsilon_L^{\alpha-1}}{\sigma_0 \alpha - 1} \quad (44)$$

and

$$\pi(\mu, \sigma | \alpha) = \frac{1}{\Delta\mu \sigma_0} \frac{\alpha - 1}{\varepsilon_H^{1-\alpha} - \varepsilon_L^{\alpha-1}} \left(\frac{\sigma}{\sigma_0} \right)^{\alpha-2}. \quad (45)$$

Notice that in the special case of $\alpha=1$, the prior over σ reduces to $d\sigma/\sigma$ which is called the Jeffreys prior and is usually introduced by the requirement of invariance under scale transformations, $\sigma \rightarrow \lambda\sigma$.

Writing $\varepsilon \stackrel{\text{def}}{=} (\varepsilon_L \varepsilon_H)^{1/2}$, the prior for α can be obtained from Eq. (32),

$$\gamma(\alpha) = \frac{1}{(\alpha - 1)^2} - \left(\frac{2 \log \varepsilon}{\varepsilon^{1-\alpha} - \varepsilon^{\alpha-1}} \right)^2 \quad (46)$$

and from Eqs. (26) and (35),

$$\pi(\alpha) = \frac{\gamma^{1/2}(\alpha) \varepsilon^{1-\alpha} - \varepsilon^{\alpha-1}}{z \alpha - 1} \exp \left[\frac{1}{\alpha - 1} + \alpha \frac{\varepsilon^{1-\alpha} + \varepsilon^{\alpha-1}}{\varepsilon^{1-\alpha} - \varepsilon^{\alpha-1}} \log \varepsilon \right], \quad (47)$$

where the normalization z has been suitably redefined.

Equations (46) and (47) simplify considerably when we take the limit $\varepsilon \rightarrow 0$. Clearly the same result is obtained whether we let $\varepsilon_H \rightarrow 0$ while keeping ε_L fixed, or letting $\varepsilon_L \rightarrow 0$ while keeping ε_H fixed, or even allowing $\varepsilon_H \rightarrow 0$ and $\varepsilon_L \rightarrow 0$ simultaneously. The resulting $\gamma(\alpha)$ and $\pi(\alpha)$ are

$$\gamma(\alpha) = \frac{1}{(\alpha - 1)^2}, \quad (48)$$

and

$$\pi(\alpha) = \begin{cases} \frac{1}{(1 - \alpha)^2} \exp \left[\frac{1}{\alpha - 1} \right] & \text{for } \alpha < 1, \\ 0 & \text{for } \alpha \geq 1, \end{cases} \quad (49)$$

where $\pi(\alpha)$ is normalized. This is shown in Fig. 1.

$\pi(\alpha)$ reaches its maximum value at $\alpha=1/2$. Since $\pi(\alpha) \sim \alpha^{-2}$ for $\alpha \rightarrow -\infty$ the expected value of α and all higher moments diverge. This suggests that replacing the unknown α in the prior $\pi(\theta | \alpha)$ by any given numerical value $\hat{\alpha}$ is probably not a good approximation.

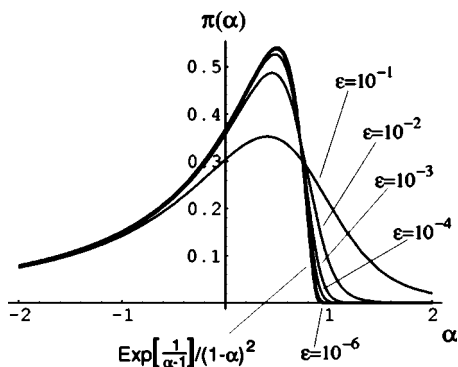


FIG. 1. The prior $\pi(\alpha)$ for various values of the cutoff parameter ε , as $\varepsilon \rightarrow 0$.

As explained in Sec. V, since α is unknown, the effective prior for $\theta=(\mu, \sigma)$ is obtained marginalizing $\pi(\mu, \sigma, \alpha)$ over α , Eq. (26). Since $\pi(\alpha)=0$ for $\alpha \geq 1$ as $\varepsilon \rightarrow 0$ we can safely take the limit $\varepsilon_H \rightarrow 0$ or $\sigma_H \rightarrow \infty$. Conversely, since $\pi(\alpha) \neq 0$ for $\alpha < 1$ we cannot take $\varepsilon_L \rightarrow 0$ or $\sigma_L \rightarrow 0$. The limit $\sigma_H \rightarrow \infty$ while keeping σ_L fixed gives

$$\pi(\mu, \sigma, \alpha) = \begin{cases} \frac{1}{\Delta\mu\sigma_L} \frac{\exp\left[\frac{1}{\alpha-1}\right]}{1-\alpha} \left(\frac{\sigma}{\sigma_L}\right)^{\alpha-2} & \text{for } \alpha < 1, \\ 0 & \text{for } \alpha \geq 1. \end{cases} \quad (50)$$

The averaged prior for μ and σ is

$$\bar{\pi}(\mu, \sigma) = \frac{1}{\Delta\mu\sigma_L} \left(\frac{\sigma_L}{\sigma}\right)^2 \int_{-\infty}^1 \frac{\exp\left[\frac{1}{\alpha-1}\right]}{1-\alpha} \left(\frac{\sigma}{\sigma_L}\right)^{\alpha} d\alpha, \quad (51)$$

which integrates to

$$\bar{\pi}(\mu, \sigma) = \frac{2}{\Delta\mu\sigma} K_0\left(2\sqrt{\log\frac{\sigma}{\sigma_L}}\right), \quad (52)$$

where K_0 is a modified Bessel function of the second kind. This is the entropic prior for the Gaussian model. The function

$$P(x) = \frac{2}{x} K_0(2\sqrt{\log x}) \quad (53)$$

is shown in Fig. 2 as a function of $x=\sigma/\sigma_L$.

$P(x)$ has an integrable singularity as $x \rightarrow 1$ where it behaves as

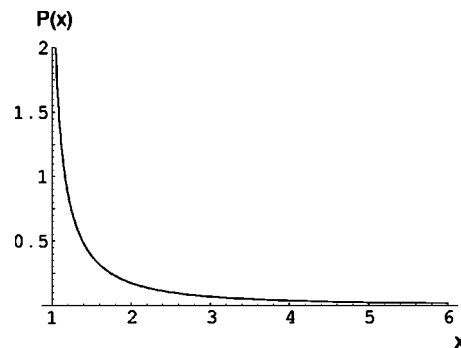


FIG. 2. The effective $\bar{\pi}(\mu, \sigma)$ is shown as $P(x) = (2/x)K_0[2\sqrt{\log(x)}]$ where $x=\sigma/\sigma_L$.

$$P(x) \approx \frac{2}{x} - \log\sqrt{\log x} - \gamma \text{ for } x \approx 1. \quad (54)$$

Since σ_L is a lower cutoff the region of large x is more relevant. The leading asymptotic behavior is given by

$$P(x) \approx \frac{\sqrt{\pi}}{x(\log x)^{1/4}} \exp -2\sqrt{\log x} \text{ for } x \gg 1. \quad (55)$$

Finally, we turn to Bayes' theorem, Eq. (25), with the prior (52) and consider using the expected values $\langle\mu\rangle$ and $\langle\sigma\rangle$ over the posterior $p(\mu, \sigma|y^N)$ as estimators for μ and σ . The integrations can be performed numerically but are not particularly illuminating. Of course, for large N the results are independent of the prior and the estimators coincide with the standard maximum likelihood results, namely $\langle\mu\rangle \rightarrow \bar{y}$ (the sample average) and $\langle\sigma\rangle^2 \rightarrow \bar{y}^2 - \bar{y}^2$ (the sample variance).

Alternatively, one can follow standard practice and marginalize the posterior $p(\mu, \sigma|y^N)$ over σ to obtain $p(\mu|y^N)$ and calculate an estimator $\hat{\mu}$ and its error bar $\Delta\mu$ from

$$\frac{d}{d\mu} \log p(\mu|y^N)|_{\hat{\mu}} = 0 \quad (56)$$

and

$$-\frac{d^2}{d\mu^2} \log p(\mu|y^N)|_{\hat{\mu}} = \frac{1}{(\Delta\mu)^2}. \quad (57)$$

The result for $\hat{\mu}$ is very simple: for any value of N the estimator $\hat{\mu}$ is the sample average, $\hat{\mu} = \bar{y}$. The result for $\Delta\mu$ is not as elegant but for large N it simplifies and reduces to $(\Delta\mu)^2 \rightarrow (\bar{y}^2 - \bar{y}^2)/N$ as expected.

VIII. FINAL REMARKS

In this paper the method of maximum relative entropy has been used to translate the information contained in the known form of the likelihood into a prior distribution for Bayesian inference. The argument follows closely the analo-

gous ME methods that have been so successful in statistical mechanics. For experiments that cannot be repeated the resulting “entropic prior” is formally identical with the Einstein fluctuation formula. For repeatable experiments, however, the expected value of the entropy of the likelihood — represented in terms of a Lagrange multiplier α — turns out to be relevant information that must be included in the analysis. As an illustration the important case of a Gaussian likelihood was treated in detail.

It may be useful to comment briefly on the differences between our entropic prior and the versions previously proposed by Skilling and by Rodríguez. Perhaps the main difference with Skilling’s prior is that, unlike ours, its use is not restricted to probability distributions but is intended for generic “positive additive distributions” including, for example, the distributions of intensities in images [15]. One problem here is that of justifying the applicability of the ME method in such a general context. Our impulse to generalize is a dangerous one; we may get away with indulging it occasionally but overindulgence will certainly lead to error. In any case, our argument in Sec. III, which consists in maximizing the entropy σ subject to a constraint $p(y, \theta) = \pi(\theta)p(y|\theta)$, makes no sense in the case of generic positive additive distributions for which there is no available product rule. A more specific problem arises from the fact that Skilling’s entropy is not, in general, dimensionless and the hyperparameter α is vaguely interpreted some sort of cutoff carrying the appropriate corrective units. Some of the difficulties, which led Skilling to seek an alternative approach, were identified in Ref. [29].

Rodríguez’s approach is closer to ours. His prior applies to probability distributions and appears to be derived from a ME principle [25]. One difference, perhaps a minor one, is his treatment of the underlying measure $m(y)$. For us $m(y)$ is not arbitrary; knowing $m(y)$ is part of knowing what data has been collected. For him $m(y)$ is just an initial guess and he suggests setting $m(y) = p(y|\theta_0)$ for some value θ_0 . The more important difference, however, is that the number of observed data n is deliberately and explicitly left unspecified. The space $\Theta \times Y^n$ over which distributions are defined, and therefore the distributions themselves, also remain unspecified. It is not clear what the maximization of an entropy over such unspecified spaces could possibly mean but a hyperparameter α is eventually introduced and it is interpreted as a “virtual number of observations supporting the initial guess θ_0 .” He proposes that α be considered as one more among the parameters θ to be inferred. As mentioned earlier this leads to the introduction of an endless chain of additional hyperparameters.

There are several directions in which the ideas of this paper can be further extended. First, we emphasize once again that the entropic priors discussed here apply to a situation where all we know about the quantities θ is that they appear as parameters in the likelihood $p(y|\theta)$, and *nothing else*. In many situations of experimental interest there exists additional relevant information beyond what is contained in the likelihood. Such information should be included as additional constraints in the maximization of the relative entropy σ in Eq. (17). The resulting modified entropic prior would

provide a better representation of our state of knowledge prior to the acquisition of the data. Indeed, the advantage of the Bayesian approach over the usual method of maximum likelihood is the possibility of including additional relevant information by replacing a flat prior by an appropriately more informative prior. There is nothing to prevent us from performing a similar improvement and going beyond the “bare” entropic priors discussed in this paper. Two kinds of additional information that are easy to include are restrictions on the range of the parameters θ and information about the known expected values of some variables $a(\theta)$. Steps in this direction were taken in Sec. V, where $a(\theta)$ is the likelihood entropy, and in Sec. VII where high and low cutoffs on the range of the Gaussian parameters were introduced.

Second, in the introduction we mentioned the interesting possibility of analyzing data y_e from different experiments, $e=1,2,\dots$, related to θ by different likelihood functions $p_e(y_e|\theta)$. Clearly this can be analyzed as a single combined experiment with likelihood $p(y_1, y_2, \dots|\theta) = p_1(y_1|\theta)p_2(y_2|\theta)\dots$ to which all our previous results apply. As we stated earlier, the mere fact that θ is measurable through one or another experiment is additional relevant information that can be taken into account.

Third, we also mentioned that problems of model selection can be tackled as an extension of the ideas described in this paper. On the basis of data y we want to select one model among several competing candidates labeled by $m=1,2,\dots$ with likelihood distributions given by $p(y|m, \theta_m)$. The answer, i.e., the probability of model m given the data y , is given by Bayes’ theorem,

$$\begin{aligned} p(m|y) &= \frac{\pi(m)}{p(y)} p(y|m) \\ &= \frac{\pi(m)}{p(y)} \int d\theta_m p(y, \theta_m|m) \\ &= \frac{1}{p(y)} \int d\theta_m \pi(m, \theta_m) p(y|m, \theta_m). \end{aligned} \quad (58)$$

This is exact. The problem is solved, at least in principle, once an entropic prior for $\pi(m, \theta_m)$ is assigned. However, the remaining practical problems associated with carrying out the actual numerical calculations could, of course, still be quite formidable.

Finally, we end with a word of caution. As in all instances of inductive inference there is the possibility that predictions based on the ME method could be wrong because not all the information relevant to the problem at hand was taken into account. This potential problem is not peculiar to the ME method; it is a problem shared by all methods of induction. Nevertheless, we are confident that there will be enormous rewards in extending the benefits of the ME method, as a method for induction, beyond its traditional territory in statistical mechanics into data analysis.

ACKNOWLEDGMENTS

Many of our comments and arguments have been inspired by Carlos C. Rodríguez, Volker Dose, and Rainer Fischer

through insightful questions and discussions which we gratefully acknowledge. A. C. also acknowledges the hospitality of the Max-Planck-Institut für Plasmaphysik during the two extended visits when most of this work was carried out.

-
- [1] For a recent review see V. Dose, Rep. Prog. Phys. **66**, 1421 (2003); for a pedagogical introduction, see D. S. Sivia, *Data Analysis, A Bayesian Tutorial* (Oxford University Press, Oxford, 1996).
- [2] *Statistical Challenges in Modern Astronomy I-III*, series edited by E. D. Feigelson and G. J. Babu (Springer, New York, 1992, 1997, 2002).
- [3] V. Dose, R. Preuss, and W. von der Linden, Phys. Rev. Lett. **81**, 3407 (1998).
- [4] G. D'Agostini, CERN Yellow Report 99-03 1999.
- [5] W. von der Linden, M. Donath, and V. Dose, Phys. Rev. Lett. **71**, 899 (1993).
- [6] R. Preuss *et al.*, Phys. Rev. Lett. **73**, 732 (1994); R. Preuss, W. Hanke, and W. von der Linden, *ibid.* **75**, 1344 (1995); R. Preuss, W. Hanke, C. Gröber, and H. G. Evertz, *ibid.* **79**, 1122 (1997).
- [7] R. Fischer, M. Mayer, W. von der Linden, and V. Dose, Phys. Rev. E **55**, 6667 (1997).
- [8] W. von der Linden, V. Dose, J. Padayachee, and V. Prozesky, Phys. Rev. E **59**, 6527 (1999); R. Fischer, K. M. Hanson, V. Dose, and W. von der Linden, *ibid.* **61**, 1152 (2000).
- [9] T. Schwarz-Selinger, R. Preuss, V. Dose, and W. von der Linden, J. Mass Spectrom. **36**, 866 (2001).
- [10] U. V. Toussaint, R. Fischer, K. Krieger, and V. Dose, New J. Phys. **1**, 11 (1999).
- [11] G. L. Bretthorst, *Bayesian Spectrum Analysis and Parameter Estimation*, Lecture Notes in Physics, Vol. 48 (Springer-Verlag, Berlin, 1988); J. Magn. Reson. (1969-1992) **88**, 552 (1990).
- [12] For a review with annotated bibliography, see, e.g., R. E. Kass and L. Wasserman, J. Am. Stat. Assoc. **91**, 1343 (1996).
- [13] J. M. Bernardo, T. Z. Irony, and N. D. Singpurwalla, J. Stat. Plan. Infer. **65**, 159 (1997).
- [14] E. T. Jaynes, IEEE Trans. Syst. Sci. Cybern. **SSC-4**, 227 (1968); J. M. Bernardo, J. R. Stat. Soc. Ser. B. Methodol. **41**, 113 (1979); A. Zellner, in *Maximum Entropy and Bayesian Methods*, edited by W. T. Grandy, Jr. and L. H. Schick (Kluwer, Dordrecht, 1991).
- [15] J. Skilling, in *Maximum Entropy and Bayesian Methods*, edited by J. Skilling (Kluwer, Dordrecht, 1989); in *Maximum Entropy and Bayesian Methods*, edited by P. F. Fougère (Kluwer, Dordrecht, 1990).
- [16] C. C. Rodríguez, in *Maximum Entropy and Bayesian Methods*, edited by J. Skilling (Kluwer, Dordrecht, 1989); in *Maximum Entropy and Bayesian Methods*, edited by P. F. Fougère (Kluwer, Dordrecht, 1990); in *Maximum Entropy and Bayesian Methods*, edited by W. T. Grandy, Jr. and L. H. Schick (Kluwer, Dordrecht, 1991); in *Maximum Entropy and Bayesian Methods*, edited by G. R. Heidbreder (Kluwer, Dordrecht, 1996).
- [17] On terminology: The terms “prior” and “posterior” are normally used in the context of Bayes’ theorem; we retain the same terminology when using ME because we are concerned with the similar goal of processing information to update from a prior to a posterior. The “method of ME” is usually understood in the restricted sense that one updates from a prior distribution that happens to be uniform. Here we adopt a broader meaning that includes updates from arbitrary priors and which involves the maximization of relative entropy. Indeed since all entropies are relative to some prior the qualifier “relative” is not needed and will henceforth be omitted.
- [18] C. E. Shannon, Bell Syst. Tech. J. **27**, 379 (1948); **27**, 623 (1948); C. E. Shannon and W. Weaver, *The Mathematical Theory of Communication* (University of Illinois Press, Urbana, 1949); N. Wiener, *Cybernetics* (MIT Press, Cambridge, 1948); L. Brillouin, *Science and Information Theory* (Academic, New York, 1956); S. Kullback, *Information Theory and Statistics* (Wiley, New York, 1959).
- [19] E. T. Jaynes, Phys. Rev. **106**, 620 (1957); E. T. Jaynes, Phys. Rev. **108**, 171 (1957); *Papers on Probability, Statistics and Statistical Physics*, edited by R. D. Rosenkrantz and E. T. Jaynes (Reidel, Dordrecht, 1983); E. T. Jaynes, *Probability Theory: The Logic of Science* (Cambridge University Press, Cambridge, 2003).
- [20] J. E. Shore and R. W. Johnson, IEEE Trans. Inf. Theory **IT26**, 26 (1980); Y. Tikochinsky, N. Z. Tishby, and R. D. Levine, Phys. Rev. Lett. **52**, 1357 (1984).
- [21] I. Csiszar, Ann. Stat. **19**, 2032 (1991).
- [22] J. Skilling, in *Maximum-Entropy and Bayesian Methods in Science and Engineering*, edited by G. J. Erickson and C. R. Smith (Kluwer, Dordrecht, 1988).
- [23] A. Caticha, in *Bayesian Inference and Maximum Entropy Methods in Science and Engineering*, AIP Conf. Proc. No. 707, edited by G. Erickson and Y. Zhai (AIP, Melville, 2004), p. 75 (e-print physics/0311093).
- [24] C. C. Rodríguez, in *Maximum Entropy and Bayesian Methods*, edited by W. von der Linden, V. Dose, R. Fischer, and R. Preuss (Kluwer, Dordrecht, 1999), Sec. 3 of “Are we cruising a hypothesis space?”
- [25] C. C. Rodríguez, in *Bayesian Inference and Maximum Entropy Methods in Science and Engineering*, AIP Conf. Proc. No. 617, edited by R. L. Fry (AIP, Melville, 2002), p. 410 (e-print physics/0201016).
- [26] A. Caticha, in *Bayesian Methods and Maximum Entropy in Science and Engineering*, AIP Conf. Proc. No. 568, edited by A. Mohammad-Djafari (AIP, Melville, 2001), p. 94 (e-print math-ph/0008017).
- [27] The number and the wording of our axioms differs from Skilling’s because we concentrate on the specific problem of ranking probability distributions while he was concerned with ranking general positive additive distributions. Proofs are given in [23].

- [28] S. Amari, *Differential-Geometrical Methods in Statistics* (Springer-Verlag, New York, 1985); for a brief derivation, see A. Caticha, in *Bayesian Methods and Maximum Entropy in Science and Engineering* (Ref. [26]), p. 72 (e-print math-ph/0008018).
- [29] J. Skilling and S. Sibisi, in *Maximum-Entropy and Bayesian Methods*, edited by K. M. Hanson and R. N. Silver (Kluwer, Dordrecht, 1996); J. Skilling, in *Maximum-Entropy and Bayesian Methods*, edited by G. J. Erickson, J. T. Ryckert, and C. R. Smith (Kluwer, Dordrecht, 1998).