

# Bayes and Discovery: Objective Bayesian Hypothesis Testing

José M. Bernardo

Universitat de València, Spain

## Abstract

Hypothesis testing is formulated from a decision theoretical viewpoint. The combined use of intrinsic discrepancy, an invariant information-based loss function, and conventional reference priors provides an objective Bayesian solution to precise hypothesis testing problems which easily integrates with the standard formulation of objective Bayesian point and region estimation.

## 1 Precise Hypothesis Testing

Let  $z$  be the available data which are assumed to have been generated as one random observation from model  $\mathcal{M}_z = \{p(z|\omega), z \in \mathcal{Z}, \omega \in \Omega\}$ . Often, but not always, data will consist of a random sample  $z = \{x_1, \dots, x_n\}$  from some distribution  $q(x|\omega)$ , with  $x \in \mathcal{X}$ ; in this case  $p(z|\omega) = \prod_{i=1}^n q(x_i|\omega)$  and  $\mathcal{Z} = \mathcal{X}^n$ . Let  $\theta(\omega)$  be the vector of interest. Without loss of generality, the model may explicitly be expressed in terms of  $\theta$  so that  $\mathcal{M}_z = \{p(z|\theta, \lambda), z \in \mathcal{Z}, \theta \in \Theta, \lambda \in \Lambda\}$ , where  $\lambda$  is some appropriately chosen nuisance parameter vector. Let  $\pi(\theta, \lambda) = \pi(\lambda|\theta)\pi(\theta)$  be the assumed prior, and let  $\pi(\theta|z)$  be the corresponding marginal posterior distribution of  $\theta$ . Appreciation of the inferential contents of  $\pi(\theta|z)$  may be enhanced by providing both point and region estimates of the vector of interest  $\theta$ , and by declaring whether or not some context-suggested specific value  $\theta_0$  is compatible with the observed data  $z$  (precise hypothesis testing). A large number of Bayesian estimation and hypothesis testing procedures have been proposed in the literature. We argue that their choice is better made in decision theoretical terms.

Let  $\ell\{\theta_0, (\theta, \lambda)\}$  describe, as a function of the (unknown) parameter values  $(\theta, \lambda)$  which have generated the available data, *the loss* to be suffered if, working with model  $\mathcal{M}_z$ , the value  $\theta_0$  were used as a proxy for the unknown value of  $\theta$ . Point estimation, region estimation and hypothesis testing procedures may all be appropriately described as specific decision problems using a common prior distribution and a common loss function of this type. The results, which are obviously all conditional on the assumed model  $\mathcal{M}_z$ , may dramatically depend on the particular choices made for both the prior and the loss functions but, given the available data  $z$ , they all only depend on those through the corresponding posterior expected loss,

$$\bar{\ell}(\theta_0|z) = \int_{\Theta} \int_{\Lambda} \ell\{\theta_0, (\theta, \lambda)\} \pi(\theta, \lambda|z) d\theta d\lambda. \quad (1)$$

As a function of  $\theta_0 \in \Theta$ , the expected loss  $\bar{\ell}(\theta_0|z)$  provides a direct measure of the relative unacceptability of all possible values of the quantity of interest in the light of the information provided by the data.

In this paper, we will concentrate on precise hypothesis testing, with objective reference priors. For a more general perspective and many examples, see Bernardo (2011) and references therein.

### 1.1 Decision Theoretic Formulation

Consider a value  $\theta_0$  of the vector of interest which deserves special consideration, either because assuming  $\theta = \theta_0$  would noticeably simplify the model, or because there are additional context-specific arguments suggesting that  $\theta = \theta_0$ . Intuitively, the value  $\theta_0$  should be judged to be *incompatible* with

the observed data  $z$  if the posterior expected loss  $\bar{\ell}(\theta_0 | z)$  of using  $\theta_0$  as a proxy for  $\theta$  is too large. This notion is now made precise.

Formally, testing the hypothesis  $H_0 \equiv \{\theta = \theta_0\}$  may be described as a decision problem where the action space  $\mathcal{A} = \{a_0, a_1\}$  contains only two elements: to accept ( $a_0$ ) or to reject ( $a_1$ ) the hypothesis under scrutiny. Foundations require specification of a loss function  $\ell_h\{a_i, (\theta, \lambda)\}$  measuring the consequences of accepting or rejecting  $H_0$  as a function of the actual parameter values. By assumption,  $a_0$  means to *act as if*  $H_0$  were true, that is to work with the model  $\mathcal{M}_0 = \{p(z | \theta_0, \lambda_0), z \in \mathcal{Z}, \lambda_0 \in \Lambda\}$ , while  $a_1$  means to reject this simplification and to keep working with model  $\mathcal{M}_z = \{p(z | \theta, \lambda), z \in \mathcal{Z}, \theta \in \Theta, \lambda \in \Lambda\}$ . Alternatively, an already established model  $\mathcal{M}_0$  may have been embedded into a more general model  $\mathcal{M}_z$ , constructed to include promising departures from  $\theta = \theta_0$ , and it is required to verify whether presently available data  $z$  are still compatible with  $\theta = \theta_0$ , or whether the extension to  $\theta \in \Theta$  is really necessary. Given the available data  $z$ , the optimal action will be to reject the hypothesis considered if (and only if) the expected posterior loss of accepting ( $a_0$ ) is larger than that of rejecting ( $a_1$ ), so that  $\int_{\Theta} \int_{\Lambda} [\ell_h\{a_0, (\theta, \lambda)\} - \ell_h\{a_1, (\theta, \lambda)\}] \pi(\theta, \lambda | z) d\theta d\lambda > 0$ . Hence, only the loss difference  $\Delta\ell_h\{\theta_0, (\theta, \lambda)\} = \ell_h\{a_0, (\theta, \lambda)\} - \ell_h\{a_1, (\theta, \lambda)\}$ , which measures the *advantage* of rejecting  $H_0 \equiv \{\theta = \theta_0\}$  as a function of the parameter values, must be specified. The hypothesis  $H_0$  should be rejected whenever the expected advantage of rejecting is positive. Without loss of generality, the function  $\Delta\ell_h$  may be written in the form

$$\Delta\ell_h\{\theta_0, (\theta, \lambda)\} = \ell\{\theta_0, (\theta, \lambda)\} - \ell_0$$

where, as mentioned above (and precisely as in estimation problems),  $\ell\{\theta_0, (\theta, \lambda)\}$  describes the non-negative loss to be suffered if  $\theta_0$  were used as a proxy for  $\theta$ . Since  $\ell\{\theta_0, (\theta_0, \lambda)\} = 0$ , so that  $\Delta\ell_h\{\theta_0, (\theta_0, \lambda)\} = -\ell_0$ , the value  $\ell_0 > 0$  describes (in the same loss units) the context-dependent non-negative advantage of accepting  $\theta = \theta_0$  when it is true. With this formulation, the optimal action is to reject  $\theta = \theta_0$  whenever the expected value of  $\ell\{\theta_0, (\theta, \lambda)\} - \ell_0$  is positive, i.e., whenever  $\bar{\ell}(\theta_0 | z)$ , the posterior expectation of  $\ell\{\theta_0, (\theta, \lambda)\}$ , is larger than  $\ell_0$ . Thus, as intuition suggested, the solution to the precise hypothesis testing decision problem posed is found in terms of the value of the expected loss  $\bar{\ell}(\theta_0 | z)$  of using  $\theta_0$  as a proxy for the unknown value of  $\theta$ .

Using the zero-one loss function,  $\ell\{\theta_0, (\theta, \lambda)\} = 0$  if  $\theta = \theta_0$ , and  $\ell\{\theta_0, (\theta, \lambda)\} = 1$  otherwise, so that the loss advantage of rejecting  $\theta_0$  is a constant whenever  $\theta \neq \theta_0$  and zero otherwise, leads to rejecting  $H_0$  if (and only if)  $\Pr(\theta = \theta_0 | z) < p_0$  for some context-dependent  $p_0$ . Notice that, using this particular loss function, if one is to avoid a systematic rejection of  $H_0$  (whatever the data), the prior probability  $\Pr(\theta = \theta_0)$  must be *strictly positive*. If  $\theta$  is a continuous parameter this requires the use of a non-regular “sharp” prior, concentrating a positive probability mass at  $\theta_0$ . With no mention of the (rather naïve) loss structure which is implicit in the formulation, this type of solution was early advocated by Jeffreys (1961). Notice however, that this formulation implies the use of radically different priors for hypothesis testing than those used for estimation, and a different prior for each value to be tested. Moreover, this formulation is known to lead to the difficulties associated with Lindley’s paradox (Lindley, 1957; Bartlett, 1957; Robert, 1993).

There are many real world situations where there is really a concentration of prior probability around particular value, and a sound Bayesian analysis should then certainly use this information. Under some conditions, those situations may well be described with a probability mass at a (measure zero) point. However, even in these cases, robustness concerns suggest that it may well be worth exploring the consequences of using a regular reference prior with the same data, if only to verify the possible dependence of the conclusions reached on the particular prior assumptions made.

Using the quadratic loss function leads to rejecting a  $\theta_0$  value whenever its Euclidean distance to  $E[\theta | z]$ , the posterior expectation of  $\theta$ , is sufficiently large. Observe that the use of continuous loss functions (such as the quadratic loss) permits the use in hypothesis testing of precisely the same priors

that are used in estimation, and the same prior for all values to be tested. In general, the Bayes test criterion is not invariant under one-to-one transformations. Thus, if  $\phi(\theta)$  is a one-to-one transformation of  $\theta$ , rejecting  $\theta = \theta_0$  does not generally imply rejecting  $\phi(\theta) = \phi(\theta_0)$ . Once more, invariant Bayes test procedures are available by using invariant loss functions.

The threshold constant  $\ell_0$ , which is used to decide whether or not an expected loss is too large, is part of the specification of the decision problem, and should be context-dependent. However, as shown below, a judicious choice of the loss function leads to calibrated expected losses, where the relevant threshold constant has an immediate, operational interpretation.

## 2 The Intrinsic Divergence Loss

Conditional on model  $\mathcal{M}_z = \{p(z | \theta, \lambda), z \in \mathcal{Z}, \theta \in \Theta, \lambda \in \Lambda\}$ , the required loss function  $\ell\{\theta_0, (\theta, \lambda)\}$  should describe, in terms of the unknown parameter values  $(\theta, \lambda)$  which have generated the available data, the loss to be suffered if, working with model  $\mathcal{M}_z$ , the value  $\theta_0$  were used as a proxy for  $\theta$ . It may naïvely appear that what is needed is just some measure of the discrepancy between  $\theta_0$  and  $\theta$ . However, since all parameterizations are arbitrary, what is really required is some measure of the discrepancy between the *models* labelled by  $\theta$  and by  $\theta_0$ . By construction, such a discrepancy measure will be independent of the particular parameterization used. Robert (1996) coined the word *intrinsic* to refer to those model-based loss functions. They are always invariant under one-to-one reparameterizations.

A reasonable measure of the dissimilarity  $\delta\{p_z, q_z\}$  between two probability densities  $p(z)$  and  $q(z)$  for a random vector  $z \in \mathcal{Z}$  should surely be non-negative, zero if (and only if),  $p(z) = q(z)$  almost everywhere, and preferably symmetric. Moreover it should be invariant under one-to-one transformations of  $z$ ; indeed, if  $y = \mathbf{y}(z)$  is such a transformation and  $J$  is the appropriate Jacobian,  $p_y = p_z/|J|$ , and  $q_y = q_z/|J|$  are expressions of precisely the same uncertainties and, therefore, one should certainly have  $\delta\{p_z, q_z\} = \delta\{p_y, q_y\}$ . Finally, it should also be possible to use  $\delta$  to compare densities with strictly nested supports, since many approximations are precisely obtained by restricting the original support to some strict subspace. These desiderata are all satisfied by the *intrinsic discrepancy* (Bernardo and Rueda, 2002), a divergence measure which has both an information theoretical justification, and a simple operational interpretation in terms of average log-density ratios.

**Definition 1** *The intrinsic discrepancy  $\delta\{p_1, p_2\}$  between two probability distributions for the random vector  $z$  with densities  $p_1(z)$ ,  $z \in \mathcal{Z}_1$ , and  $p_2(z)$ ,  $z \in \mathcal{Z}_2$ , is*

$$\delta\{p_1, p_2\} = \min [\kappa\{p_1 | p_2\}, \kappa\{p_2 | p_1\}] \quad (2)$$

where  $\kappa\{p_j | p_i\} = \int_{\mathcal{Z}_i} p_i(z) \log[p_i(z)/p_j(z)] dz$  is the Kullback-Leibler (KL) directed logarithmic divergence of  $p_j$  from  $p_i$ . The intrinsic discrepancy between a probability distribution  $p$  and a family of distributions  $\mathcal{F} = \{q_i, i \in I\}$  is the intrinsic discrepancy between  $p$  and the closest of them,

$$\delta\{p, \mathcal{F}\} = \inf_{q \in \mathcal{F}} \delta\{p, q\}.$$

The intrinsic discrepancy  $\delta\{p_1, p_2\}$  is the minimum average log density ratio of one density over the other, and has an operative interpretation as the minimum amount of information (in natural information units or *nits*) expected to be required to discriminate between  $p_1$  and  $p_2$ . This may be used to define an appropriate loss function for the decision problem considered in this paper as the intrinsic discrepancy between the model, labelled by  $(\theta, \lambda)$ , and the family  $\mathcal{M}_0$  of models which satisfy the hypothesis to be tested:

**Definition 2** *Consider  $\mathcal{M}_z = \{p(z | \theta, \lambda), z \in \mathcal{Z}, \theta \in \Theta, \lambda \in \Lambda\}$ . The intrinsic discrepancy loss of using  $\theta_0$  as a proxy for  $\theta$  is the intrinsic discrepancy between the true model and the class of models*

with  $\theta = \theta_0$ ,  $\mathcal{M}_0 = \{p(z | \theta_0, \lambda_0), z \in \mathcal{Z}, \lambda_0 \in \Lambda\}$ ,

$$\ell_\delta\{\theta_0, (\theta, \lambda) | \mathcal{M}_z\} = \delta\{p_z(\cdot | \theta, \lambda), \mathcal{M}_0\} = \inf_{\lambda_0 \in \Lambda} \delta\{p_z(\cdot | \theta_0, \lambda_0), p_z(\cdot | \theta, \lambda)\}. \quad (3)$$

Notice the complete generality of Definition 2; this may be used with either discrete or continuous data models (in the discrete case, the integrals in Definition 1 will obviously be sums), and with either discrete or continuous parameter spaces of any dimensionality.

The intrinsic discrepancy loss has many attractive invariance properties. For any one-to-one reparameterization of the form  $\phi = \phi(\theta)$  and  $\psi = \psi(\theta, \lambda)$ ,

$$\ell_\delta\{\theta_0, (\theta, \lambda) | \mathcal{M}_z\} = \ell_\delta\{\phi_0, (\phi, \psi) | \mathcal{M}_z\},$$

so that the use of this loss function will lead to estimation and hypothesis testing procedures which are *invariant* under those transformations. Moreover, if  $t = t(z)$  is a sufficient statistic for model  $\mathcal{M}_z$ , one may equivalently work with the marginal model  $\mathcal{M}_t = \{p(t | \theta, \lambda), t \in \mathcal{T}, \theta \in \Theta, \lambda \in \Lambda\}$  since, in that case,

$$\ell_\delta\{\theta_0, (\theta, \lambda) | \mathcal{M}_z\} = \ell_\delta\{\theta_0, (\theta, \lambda) | \mathcal{M}_t\}.$$

Computations are often simplified by using the additive property of the intrinsic discrepancy loss : if data consist of a random sample  $z = \{x_1, \dots, x_n\}$  from some underlying model  $\mathcal{M}_x$ , so that  $\mathcal{Z} = \mathcal{X}^n$ , and  $p(z | \theta, \lambda) = \prod_{i=1}^n p(x_i | \theta, \lambda)$ , then

$$\ell_\delta\{\theta_0, (\theta, \lambda) | \mathcal{M}_z\} = n \ell_\delta\{\theta_0, (\theta, \lambda) | \mathcal{M}_x\}.$$

An interesting interpretation of the intrinsic discrepancy loss follows directly from Definitions 1 and 2. Indeed,  $\ell_\delta\{\theta_0, (\theta, \lambda) | \mathcal{M}_z\}$  is just the minimum log-likelihood ratio which may be expected under repeated sampling between the true model, identified by  $(\theta, \lambda)$ , and the class of models which have  $\theta = \theta_0$ . Thus, *the intrinsic discrepancy loss formalizes the use of the minimum average log-likelihood ratio under sampling as a general loss function*.

In particular, a suggested value  $\theta_0$  for the vector of interest should be judged to be incompatible with the observed data  $z$  if  $\bar{\ell}_\delta(\theta_0 | z)$ , the posterior expectation of  $\ell_\delta\{\theta_0, (\theta, \lambda) | \mathcal{M}_z\}$ , is larger than a suitably chosen constant  $\ell_0$ . For instance, if for some arbitrary  $k$ ,  $\ell_0 = \log[10^k]$ , then  $\theta_0$  would be rejected whenever, given the observed data, the minimum sampling average likelihood ratio against  $\theta = \theta_0$ , may be expected to be larger than about  $10^k$ . Conventional choices for  $\ell_0$  are  $\{\log 100, \log 1000, \log 10000\} \approx \{4.6, 6.9, 9.2\}$ .

Under regularity conditions, the intrinsic discrepancy loss has an alternative expression which is generally much simpler to compute:

**Theorem 1** (Juárez, 2004, Sec. 2.4) *If the support of  $p(z | \theta, \lambda)$  is convex for all  $(\theta, \lambda)$ , then the intrinsic discrepancy loss may also be written as*

$$\ell_\delta\{\theta_0, (\theta, \lambda) | \mathcal{M}_z\} = \min \left[ \inf_{\lambda_0 \in \Lambda} \kappa\{\theta_0, \lambda_0 | \theta, \lambda\}, \inf_{\lambda_0 \in \Lambda} \kappa\{\theta, \lambda | \theta_0, \lambda_0\} \right], \quad (4)$$

where  $\kappa\{\theta_j, \lambda_j | \theta_i, \lambda_i\}$  is the KL-divergence of  $p_z(\cdot | \theta_j, \lambda_j)$  from  $p_z(\cdot | \theta_i, \lambda_i)$ .

When there is no danger of confusion,  $\mathcal{M}_z$  may be dropped from the notation and  $\ell_\delta\{\theta_0, (\theta, \lambda) | \mathcal{M}_z\}$  may be written  $\ell_\delta\{\theta_0, (\theta, \lambda)\}$ , but the dependence on the model of intrinsic losses should always be kept in mind.

In the important case of a multivariate normal model with known covariance matrix, the intrinsic discrepancy loss is proportional to the Mahalanobis distance:

**Example 1 (Multivariate normal model).** Let  $\mathbf{z} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  be a random sample from a  $k$ -variate normal distribution  $N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$  with known covariance matrix  $\boldsymbol{\Sigma}$ . The KL divergence of  $N(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma})$  from  $N(\mathbf{x} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma})$  is  $\kappa\{\boldsymbol{\mu}_j | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}\} = \frac{1}{2}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)^t \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)$ . Since this is symmetric, and the intrinsic discrepancy is additive,

$$\delta\{\boldsymbol{\mu}_0, \boldsymbol{\mu} | \boldsymbol{\Sigma}\} = \frac{n}{2}(\boldsymbol{\mu}_0 - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_0 - \boldsymbol{\mu}),$$

which is  $n/2$  times the Mahalanobis distance between  $\boldsymbol{\mu}_0$  and  $\boldsymbol{\mu}$ .

## 2.1 Approximations

Under regularity conditions, the result of Example 1 may be combined with conventional asymptotic results to obtain large sample approximations to intrinsic discrepancy losses:

**Theorem 2** (Bernardo, 2011) *Let data  $\mathbf{z} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  consist of a random sample from  $p(\mathbf{x} | \boldsymbol{\theta}, \boldsymbol{\lambda})$ , let  $F(\boldsymbol{\theta}, \boldsymbol{\lambda})$  be the corresponding Fisher matrix, and let  $V(\boldsymbol{\theta}, \boldsymbol{\lambda}) = F^{-1}(\boldsymbol{\theta}, \boldsymbol{\lambda})$  be its inverse. Then, for large  $n$  and under conditions for asymptotic normality,*

$$\ell\{\boldsymbol{\theta}_0, (\boldsymbol{\theta}, \boldsymbol{\lambda}) | \mathcal{M}_{\mathbf{z}}\} \approx \frac{n}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)^t V_{\boldsymbol{\theta}\boldsymbol{\theta}}^{-1}(\boldsymbol{\theta}, \boldsymbol{\lambda})(\boldsymbol{\theta} - \boldsymbol{\theta}_0),$$

where  $V_{\boldsymbol{\theta}\boldsymbol{\theta}}$  is the submatrix of  $V(\boldsymbol{\theta}, \boldsymbol{\lambda})$  which corresponds to the vector of interest  $\boldsymbol{\theta}$ .

The invariance of the intrinsic discrepancy loss under reparameterization may be exploited to improve the approximation above, by simply choosing a parameterization where the asymptotic convergence to normality is faster. The following result is a one-parameter example of this technique, which makes use of the variance stabilization transformation.

**Theorem 3** (Bernardo, 2005b) *Let  $\mathbf{z} = \{x_1, \dots, x_n\}$  be a random sample of size  $n$  from model  $p(x | \theta)$ , and let  $\hat{\theta}_n = \hat{\theta}_n(\mathbf{z})$  be an asymptotically sufficient consistent estimator of  $\theta$ , whose sampling distribution is asymptotically normal with standard deviation  $s(\theta)/\sqrt{n}$ . Define  $\phi(\theta) = \int^\theta s(y)^{-1} dy$ . Then,*

$$\ell\{\hat{\theta}_n, \theta | \mathcal{M}_{\mathbf{z}}\} = \frac{n}{2} [\phi(\hat{\theta}_n) - \phi(\theta)]^2 + o_p(1).$$

## 3 Reference Analysis

Foundations indicate that the prior distribution should describe available prior knowledge. In many situations however, either the available prior information on the quantity of interest is too vague or too complex to warrant the effort required to formalize it, or it is too subjective to be useful in scientific communication. An ‘‘objective’’ procedure is therefore often required, where the prior function is intended to describe a situation where there is no relevant information about the quantity of interest. Objectivity is an emotionally charged word, and it should be explicitly qualified whenever it is used. No statistical analysis is really objective, since both the experimental design and the model assumed have very strong subjective inputs. However, frequentist procedures are often branded as ‘‘objective’’ just because their conclusions are only conditional on the model assumed and the data obtained. Bayesian methods where the prior function is directly derived from the assumed model are objective in this limited, but precise sense. For lively discussions of this, and related issues, see Bernardo (1997), Berger (2006), and ensuing discussions.

There is a vast literature devoted to the formulation of objective priors; relevant pointers are included in Bernardo and Smith (1994, Sec. 5.6), Kass and Wasserman (1996), Datta and Mukerjee (2004), Bernardo (2005a), Berger (2006), Ghosh, Delampady and Samanta (2006), and references therein. Reference analysis, introduced by Bernardo (1979) and further developed by Berger and Bernardo (1989,

1992a,b,c), Sun and Berger (1998) and Berger, Bernardo and Sun (2009, 2011a,b), has been one of the most popular approaches for developing objective priors.

We will not repeat here arguments for reference analysis, but it may be worth synthesizing the basic definition and briefly reviewing some recent developments.

Note first that the same mathematical concepts which lie behind the definition of the intrinsic discrepancy provide the intuitive basis for the definition of reference priors. Indeed, for the one parameter model  $\mathcal{M} = \{p(\mathbf{z} | \theta), \mathbf{z} \in \mathcal{Z}, \theta \in \Theta \subset \mathfrak{R}\}$ , the intrinsic discrepancy  $I\{p_\theta | \mathcal{M}\} = \delta\{p(\mathbf{z}, \theta), p(\mathbf{z})p(\theta)\}$  between the joint prior  $p(\mathbf{z}, \theta)$  and the product of their marginals  $p(\mathbf{z})p(\theta)$  is a functional of the prior  $p(\theta)$  which measures the association between the data and the parameter and hence, the amount of information that, given prior  $p(\theta)$ , data  $\mathbf{z}$  may be expected to provide about  $\theta$ . If one considers  $k$  independent observations from  $\mathcal{M}$  then, as  $k$  increases,  $I\{p_\theta | \mathcal{M}^k\}$  will approach the *missing information* about  $\theta$  which repeated sampling from  $\mathcal{M}$  could provide. If  $\pi_k(\theta)$  denotes the prior which maximizes  $I\{p_\theta | \mathcal{M}^k\}$ , the sequence  $\{\pi_k(\theta)\}_{k=1}^\infty$  will converge to that prior function which maximizes the missing information about  $\theta$ , and this is defined to be the reference prior  $\pi(\theta | \mathcal{M})$ .

**Theorem 4** (Berger, Bernardo and Sun, 2009). *Let  $\mathbf{z}^{(k)} = \{z_1, \dots, z_k\}$  denote  $k$  conditionally independent observations from  $\mathcal{M}_z$ . Then, the reference prior is defined as an appropriate limit of*

$$\pi_k(\theta) \propto \exp \left\{ \mathbb{E}_{\mathbf{z}^{(k)} | \theta} [\log p_h(\theta | \mathbf{z}^{(k)})] \right\} \quad (5)$$

where  $p_h(\theta | \mathbf{z}^{(k)}) \propto \prod_{i=1}^k p(z_i | \theta) h(\theta)$  is the posterior which corresponds to any arbitrarily chosen prior function  $h(\theta)$  which makes the posterior proper for any  $\mathbf{z}^{(k)}$ .

Theorem 4 implies that the reference prior at a particular point  $\theta$  is proportional to the *logarithmic average* of the posterior density which this point would have under repeated sampling, if this  $\theta$  value were the true parameter value. The parameter values which could be expected to get relatively large asymptotic posterior densities if they were true, will then precisely be those with relatively large reference prior densities.

The result in Theorem 4 makes very simple the numerical derivation of a one-parameter reference prior. One first chooses some formal prior  $h(\theta)$ , maybe one for which exact or approximate posterior computation is easy, and a relatively large number of replications  $k$ . For each particular  $\theta$  value whose reference prior is desired, one generates a collection  $\{\mathbf{z}_1^{(k)}, \dots, \mathbf{z}_s^{(k)}\}$  of  $s$  replications  $\mathbf{z}_i^{(k)} = \{z_{i1}, \dots, z_{ik}\}$  of size  $k$  from the original model  $p(\mathbf{z} | \theta)$ , computes the corresponding  $s$  posterior densities at  $\theta$ ,  $\{p_h(\theta | \mathbf{z}_j^{(k)})\}_{j=1}^s$ , and approximates the reference prior at this point by its logarithmic average,

$$\pi(\theta) \approx \exp \left\{ \frac{1}{s} \sum_{j=1}^s \log p_h(\theta | \mathbf{z}_j^{(k)}) \right\}. \quad (6)$$

Under regularity conditions explicit formulae for the reference priors are readily available. In particular, if the posterior distribution of  $\theta$  given a random sample of size  $n$  from  $p(\mathbf{x} | \theta)$  is asymptotically normal with standard deviation  $s(\tilde{\theta}_n)/\sqrt{n}$ , where  $\tilde{\theta}_n$  is a consistent estimator of  $\theta$ , then the reference prior is  $\pi(\theta) = s(\theta)^{-1}$ . This includes as a particular case the famous Jeffreys-Perks prior (Jeffreys, 1946, independently formulated by Perks, 1947)

$$\pi(\theta) \propto i(\theta)^{1/2}, \quad i(\theta) = \mathbb{E}_{\mathbf{x} | \theta} [-\partial^2 \log p(\mathbf{z} | \theta) / \partial \theta^2]. \quad (7)$$

Similarly, if  $p(\mathbf{x} | \theta)$  is a non regular model with a support  $S(\theta)$  which depends on the parameter in the form  $S(\theta) = \{x; a_1(\theta) < x < a_2(\theta)\}$ , where the  $a_i(\theta)$ 's are monotone functions of  $\theta$  and  $S(\theta)$  is either increasing or decreasing then, under regularity conditions (Ghosal and Samanta, 1997), the reference prior is

$$\pi(\theta) \propto \mathbb{E}_{\mathbf{x} | \theta} [|\partial \log p(\mathbf{z} | \theta) / \partial \theta|]. \quad (8)$$

In multiparameter problems, reference priors depend of the quantity of interest, a necessary feature in the construction of objective priors, if one is to prevent unacceptable behaviour in the posterior, such as marginalization paradoxes (Dawid, Stone and Zidek, 1973) or strong inconsistencies (Stone, 1976).

If the model has more than one parameter, the required joint reference prior is derived sequentially. Thus, if the model is  $p(z | \theta, \lambda)$  and  $\theta$  is the quantity of interest, one works conditionally on  $\theta$  and uses the one-parameter algorithm to derive the *conditional reference prior*  $\pi(\lambda | \theta)$ . If this is proper, it is used to obtain the *integrated model*  $p(z | \theta) = \int_{\Lambda} p(z | \theta, \lambda) \pi(\lambda | \theta) d\lambda$ , to which the one-parameter algorithm is applied again to obtain the *marginal reference prior*  $\pi(\theta)$ . The *joint reference prior* to compute the reference posterior for  $\theta$  is then defined to be  $\pi(\lambda | \theta) \pi(\theta)$ . If  $\pi(\lambda | \theta)$  is not proper, one proceeds similarly within a compact approximation to the parameter space (where all reference priors will be proper) and then derives the corresponding limiting result.

In general, reference priors are sequentially derived with respect to an ordered parameterization. Thus, given a model  $\mathcal{M}_{\mathbf{z}} = \{p(z | \omega), z \in \mathcal{Z}, \omega \in \Omega\}$  with  $m$  parameters, the reference prior with respect to a particular ordered parameterization  $\phi(\omega) = \{\phi_1, \dots, \phi_m\}$  (where the  $\phi_i$ 's are ordered by inferential importance) is sequentially obtained as  $\pi(\phi) = \pi(\phi_m | \phi_{m-1}, \dots, \phi_1) \times \dots \times \pi(\phi_2 | \phi_1) \pi(\phi_1)$ . Unless all reference priors turn out to be proper, the model must be endowed with an appropriate compact approximation to the parameter space  $\{\Omega_j\}_{j=1}^{\infty} \subset \Omega$ , which should remain the same for all reference priors obtained within the same model. Berger and Bernardo (1992c) describe the relevant algorithm for regular multiparameter models where asymptotic normality may be established. In typical applications,  $\theta = \phi_1$  will be the quantity of interest, and the joint reference prior  $\pi(\phi)$ , which is often denoted  $\pi_{\theta}(\phi)$  to emphasize the role of  $\theta$ , is a just a technical device to produce the desired one-dimensional marginal reference posterior  $\pi(\theta | z)$  of the quantity of interest.

#### 4 Objective Bayesian Hypothesis Testing

With the loss function chosen to be the intrinsic discrepancy loss, all that is required to define an objective Bayesian testing procedure is to specify an objective prior distribution. It will not come as a surprise that we recommend the use of a reference prior. Thus, one must obtain the posterior expectation of the intrinsic discrepancy loss with respect to the appropriate joint reference posterior

$$d(\theta_0 | z) = \int_{\Theta} \int_{\Lambda} \ell_{\delta}\{\theta_0, (\theta, \lambda) | \mathcal{M}_{\mathbf{z}}\} \pi(\theta, \lambda | z) d\theta d\lambda. \quad (9)$$

and decide whether or not this is big enough to reject that  $\theta = \theta_0$ . The function  $d(\theta_0 | z)$  is the relevant *intrinsic* test statistic, a direct measure of the incompatibility of  $\theta_0$  with the data  $z$  in terms of the expected average log-likelihood ratio against the null.

In one parameter problems, the reference prior is unique and the solution is therefore conceptually immediate. The following toy example is intended to illustrate the general procedure:

**Example 2 (Poisson data).** Let  $z = \{x_1, \dots, x_n\}$  be a random sample from a Poisson model, so that  $p(x | \lambda) = \text{Po}(x | \lambda) = e^{-\lambda} \lambda^x / x!$ . This is a regular model, and using (7), the reference prior is immediately found to be  $\pi(\lambda) = \lambda^{-1/2}$ . This leads to the gamma reference posterior  $\pi(\lambda | z) = \pi(\lambda | t, n) = \text{Ga}(\lambda | t + 1/2, n) \propto e^{-n\lambda} \lambda^{t-1/2}$ , with  $t = \sum_{j=1}^n x_j$ .

Using Definition 2 and the additive property of the intrinsic discrepancy, the intrinsic discrepancy loss of using  $\lambda_0$  as a proxy for  $\lambda$  with a random sample of size  $n$  from a Poisson distribution with parameter  $\lambda$  is

$$\delta\{\lambda_0, \lambda | \mathcal{M}_{\mathbf{z}}\} = n \delta\{\lambda_0, \lambda | \mathcal{M}_x\} = n \min \left\{ E_{x | \lambda} \left[ \log \frac{\text{Po}(x | \lambda)}{\text{Po}(x | \lambda_0)} \right], E_{x | \lambda_0} \left[ \log \frac{\text{Po}(x | \lambda_0)}{\text{Po}(x | \lambda)} \right] \right\}$$

which immediately yields

$$\delta\{\lambda_0, \lambda | \mathcal{M}_{\mathbf{z}}\} = \begin{cases} n(\lambda - \lambda_0 + \lambda_0 \log \frac{\lambda_0}{\lambda}) & \text{if } \lambda_0 \leq \lambda, \\ n(\lambda_0 - \lambda + \lambda \log \frac{\lambda}{\lambda_0}) & \text{if } \lambda_0 \geq \lambda. \end{cases} \quad (10)$$

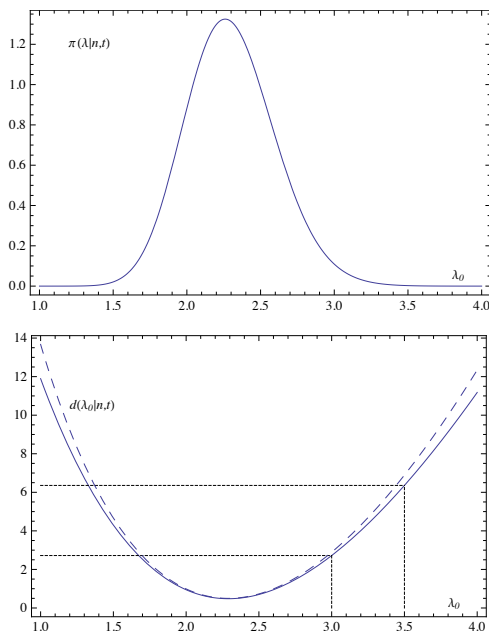


Figure 1: Posterior reference analysis of the parameter of a Poisson model.

a non-negative concave function of  $\lambda$  and  $\lambda_0$ , with minimum equal to zero when  $\lambda = \lambda_0$ .

The intrinsic statistic  $d(\lambda_0 | z)$  is the corresponding reference posterior expectation,

$$d(\lambda_0 | z) = \int_0^\infty \delta\{\lambda_0, \lambda | \mathcal{M}_z\} \text{Ga}(\lambda | t + 1/2, n) d\lambda. \quad (11)$$

This has no simple analytical expression, but is easily computed by numerical integration. Moreover, using Theorem 3 and the fact that the sampling distribution of the sufficient and consistent mle,  $\hat{\lambda} = \bar{x}$  is asymptotically normal with standard deviation  $\sqrt{\lambda}/\sqrt{n}$ , one finds

$$\begin{aligned} d(\lambda_0 | z) &\approx \int_0^\infty \frac{n}{2} (2\sqrt{\lambda_0} - 2\sqrt{\lambda})^2 \text{Ga}(\lambda | t + 1/2, n) d\lambda \\ &= 1 + 2t + 2n\lambda_0 - 4\sqrt{n\lambda_0} \frac{\Gamma(t+1)}{\Gamma(t+1/2)} \end{aligned} \quad (12)$$

$$\approx 1 + 2t + 2n\lambda_0 - 4\sqrt{n\lambda_0} \left( \sqrt{t} + \frac{1}{8\sqrt{t}} \right). \quad (13)$$

To illustrate the type of results obtained, a sample of size  $n = 25$  was simulated from a Poisson distribution with parameter  $\lambda = 2$  resulting in  $t = 57$ . The corresponding reference posterior density is plotted in the top panel of Figure 1. The expected intrinsic discrepancy loss  $d(\lambda_0 | t, n)$  (both computed from (11) (continuous line) and analytically approximated with (13) (dashed line)) are plotted in bottom panel of Figure 1. It may be appreciated that, even with this rather small sample size, the approximation is quite good.

Suppose that the value  $\lambda_0 = 3$  is to be tested. The corresponding intrinsic statistic is  $d(3 | z) = 2.72 \approx \log(15)$ ; thus the average likelihood ratio against the  $H_0 \equiv \{\lambda = 3\}$  may be expected to be about 15, not really strong evidence against this value, even if this may be seen to be in the right tail of the reference posterior of  $\lambda_0$ . On the other hand, if the value to test is  $\lambda_0 = 3.5$ , the corresponding intrinsic statistic is  $d(3.5 | z) = 6.36 \approx \log(576)$ ; hence the average likelihood ratio against this value may be expected to be about 576 and the hypothesis  $H_0 \equiv \{\lambda = 3.5\}$  should be rejected in most scenarios. Notice that, in marked difference to conventional testing using Bayes factors, the same prior has been used to test these two possible parameter values (as it would obviously be for any other value).

The following example illustrates the use of the methods described to derive a new solution to a classical precise testing problem.

**Example 3 (Equality of Normal means).** Let  $z = \{\mathbf{x}, \mathbf{y}\}$  be two independent random samples,  $\mathbf{x} = \{x_1, \dots, x_n\}$  from  $N(x | \mu_x, \sigma)$ , and  $\mathbf{y} = \{y_1, \dots, y_m\}$  from  $N(x | \mu_y, \sigma)$ , and suppose that one is interested in

comparing the two means. In particular, one may be interested in testing whether or not the precise hypothesis  $H_0 \equiv \{\mu_x = \mu_y\}$  is compatible with available data  $\mathbf{z}$ . Using the additive property of the intrinsic discrepancy loss and the fact that the KL divergence between two normals distributions with the same variance is simply  $\kappa\{\mu_j, \sigma | \mu_i, \sigma\} = (\mu_i - \mu_j)^2 / (2\sigma^2)$  to derive the logarithmic divergence of  $p(\mathbf{z} | \mu_0, \mu_0, \sigma_0)$  from  $p(\mathbf{z} | \mu_x, \mu_y, \sigma)$ , and then minimizing over both  $\mu_0$  and  $\sigma_0$  yields  $\inf_{\mu_0 \in \mathbb{R}, \sigma_0 > 0} \kappa\{\mu_0, \mu_0, \sigma_0 | \mu_x, \mu_y, \sigma\} = k_{nm} \theta^2$ , where  $k_{nm} = 2nm / (m + n)$  is the harmonic mean of the two sample sizes, and  $\theta = (\mu_x - \mu_y) / \sigma$  is the standardized difference between the two means. On the other hand,  $\inf_{\mu_0 \in \mathbb{R}, \sigma_0 > 0} \kappa\{\mu_x, \mu_y, \sigma | \mu_0, \mu_0, \sigma_0\}$  yields  $[(m + n) / 2] \log[1 + (k_{nm} / (2(m + n))) \theta^2]$ , which is always smaller. Hence, the intrinsic discrepancy loss of accepting  $H_0$  is

$$\ell_\delta\{H_0, (\mu_x, \mu_y, \sigma)\} = \ell_\delta\{H_0, \theta | \mathcal{M}\} = \frac{n + m}{2} \log \left[ 1 + \frac{k_{nm}}{2(n + m)} \theta^2 \right],$$

which reduces to  $n \log[1 + \theta^2 / 4]$  when  $n = m$ . Here, the parameter of interest is  $\theta$ . Bernardo and Pérez (2007) find that the marginal reference posterior of  $\theta$  only depends on the data through the sample sizes and  $t = t(\mathbf{z}) = (\bar{x} - \bar{y}) / (s / \sqrt{2/k_{nm}})$ , where  $s$  is the m.l.e. of  $\sigma$ . Therefore, the required marginal reference posterior of  $\theta$  is  $\pi(\theta | \mathbf{z}) = \pi(\theta | t, m, n) \propto p(t | \theta) \pi(\theta)$  where  $p(t | \theta)$  is the noncentral Student sampling distribution of  $t$ , and  $\pi(\theta) = (1 + (k_{nm} / (4(m + n))) \theta^2)^{-1/2}$  is the marginal reference prior for  $\theta$ . The posterior  $\pi(\theta | t, m, n)$  may be used to provide point and interval estimates of  $\theta$ , the standardized difference between the two means, and hence inferential statements about their relative positions.

The relevant expected loss,  $d(H_0 | t, n, m) = \int_{-\infty}^{\infty} \ell_\delta\{H_0, \theta | \mathcal{M}\} \pi(\theta | t, n, m) d\theta$ , may be used to test  $H_0$ . This has no simple analytical expression, but its value may easily be obtained by one-dimensional numerical integration. A good large sample approximation is

$$d(H_0 | t, n, m) \approx \frac{n + m}{2} \log \left[ 1 + \frac{1}{n + m} (1 + t^2) \right].$$

The sampling distribution of  $d(H_0 | t, n, m)$  is asymptotically  $(1/2)[1 + \chi_1^2(\lambda)]$ , where  $\chi_1^2(\lambda)$  is a non-central chi-squared distribution with one degree of freedom and non-centrality parameter  $\lambda = k_{nm} \theta^2 / 2$ . It follows that the expected value under sampling of  $d(H_0 | t, n, m)$  is equal to one when  $\mu_x = \mu_y$ , and increases linearly with the harmonic mean of the samples when this is not true. Thus, the testing procedure is consistent.

For many more sophisticated examples of precise hypothesis testing, see Bernardo (2011) and references therein.

## References

- Bartlett, M. (1957). A comment on D. V. Lindley's statistical paradox. *Biometrika* **44**, 533–534.
- Berger, J. O. (2006). The case for objective Bayesian analysis. *Bayesian Analysis* **1**, 385–402 and 457–464 (with discussion).
- Berger, J. O. and Bernardo, J. M. (1989). Estimating a product of means: Bayesian analysis with reference priors. *J. Amer. Statist. Assoc.* **84**, 200–207.
- Berger, J. O. and Bernardo, J. M. (1992a). Ordered group reference priors with applications to a multinomial problem. *Biometrika* **79**, 25–37.
- Berger, J. O. and Bernardo, J. M. (1992b). Reference priors in a variance components problem. *Bayesian Analysis in Statistics and Econometrics* (P. K. Goel and N. S. Yyengar, eds.) Berlin: Springer, 323–340.
- Berger, J. O. and Bernardo, J. M. (1992c). On the development of reference priors. *Bayesian Statistics 4* (J. M. Bernardo, J. O. Berger, A. P. Dawid and A. F. M. Smith, eds.) Oxford: University Press, 35–60 (with discussion).
- Berger, J. O., Bernardo, J. M. and Sun, D. (2009). The formal definition of reference priors. *Ann. Statist.* **37**, 905–938.
- Berger, J. O., Bernardo, J. M. and Sun, D. (2011a). Reference priors for discrete parameters. *J. Amer. Statist. Assoc.* (under revision).
- Berger, J. O., Bernardo, J. M. and Sun, D. (2011b). Overall reference priors. *Tech. Rep.*, Duke University, USA.

- Bernardo, J. M. (1979). Reference posterior distributions for Bayesian inference. *J. Roy. Statist. Soc. B* **41**, 113–147 (with discussion). Reprinted in *Bayesian Inference* (N. G. Polson and G. C. Tiao, eds.) Brookfield, VT: Edward Elgar, 1995, 229–263.
- Bernardo, J. M. (1997). Noninformative priors do not exist *J. Statist. Planning and Inference* **65**, 159–189 (with discussion).
- Bernardo, J. M. (2005a). Reference analysis. *Bayesian Thinking: Modeling and Computation, Handbook of Statistics* **25** (Dey, D. K. and Rao, C. R., eds.) Amsterdam: Elsevier, 17–90.
- Bernardo, J. M. (2005b). Intrinsic credible regions: An objective Bayesian approach to interval estimation. *Test* **14**, 317–384 (with discussion).
- Bernardo, J. M. (2011). Integrated objective Bayesian estimation and hypothesis testing. *Bayesian Statistics 9* (J. M. Bernardo, M. J. Bayarri, J. O. Berger, A. P. Dawid, D. Heckerman, A. F. M. Smith and M. West, eds.) Oxford University Press, 1–68, (with discussion).
- Bernardo, J. M. and Rueda, R. (2002). Bayesian hypothesis testing: A reference approach. *Internat. Statist. Rev.* **70**, 351–372.
- Bernardo, J. M. and Pérez, S. (2007). Comparing normal means: New methods for an old problem. *Bayesian Analysis* **2**, 45–58.
- Bernardo, J. M. and Smith, A. F. M. (1994). *Bayesian Theory*. Chichester: Wiley.
- Datta, G. S. and Mukerjee, R. (2004). *Probability Matching Priors: Higher Order Asymptotics*. Berlin: Springer.
- Dawid, A. P., Stone, M. and Zidek, J. V. (1973). Marginalization paradoxes in Bayesian and structural inference. *J. Roy. Statist. Soc. B* **35**, 189–233 (with discussion).
- Ghosh, J. K., Delampady, M. and Samanta, T. (2006). *An Introduction to Bayesian Analysis: Theory and Methods*. Berlin: Springer.
- Ghosal, S. and Samanta, T. (1997). Expansion of Bayes risk for entropy loss and reference prior in nonregular cases. *Statistics and Decisions* **15**, 129–140.
- Jeffreys, H. (1946). An invariant form for the prior probability in estimation problems. *Proc. Royal Soc.* **186**, 453–461.
- Jeffreys, H. (1961). *Theory of Probability* (3rd edition). Oxford: University Press.
- Juárez, M. A. (2004). *Métodos Bayesianos Objetivos de Estimación y Contraste de Hipótesis*. Ph.D. Thesis, Universitat de València, Spain.
- Kass, R. E. and Wasserman, L. (1996). The selection of prior distributions by formal rules. *J. Amer. Statist. Assoc.* **91**, 1343–1370.
- Lindley, D. V. (1957). A statistical paradox. *Biometrika* **44**, 187–192.
- Perks, W. (1947). Some observations on inverse probability, including a new indifference rule. *J. Inst. Actuaries* **73**, 285–334 (with discussion).
- Robert, C. P. (1993). A note on Jeffreys-Lindley paradox. *Statistica Sinica* **3**, 601–608.
- Robert, C. P. (1996). Intrinsic loss functions. *Theory and Decision* **40**, 192–214.
- Stone, M. (1976). Strong inconsistency from uniform priors. *J. Amer. Statist. Assoc.* **71**, 114–125 (with discussion).