

# Overall Objective Priors

JAMES O. BERGER

*Duke University, USA*

JOSE M. BERNARDO

*Universitat de València, Spain*

and

DONGCHU SUN

*University of Missouri-Columbia, USA*

## Abstract

In multi-parameter models, reference priors typically depend on the parameter or quantity of interest, and it is well known that this is necessary to produce objective posterior distributions with optimal properties. There are, however, many situations where one is simultaneously interested in all the parameters of the model or, more realistically, in several functions of them, and it would then be useful to have a single objective prior which could safely be used to produce reasonable marginal posteriors for all the quantities of interest. In this paper, we consider three methods for selecting a single objective prior and study, in a variety of problems including the multinomial problem, whether or not the resulting prior is a good approximation to the parameter-specific reference priors.

*Some key words:* Joint Reference Prior; Logarithmic Divergence; Objective Priors; Reference Analysis; Multinomial Model.

## 1 Introduction

### 1.1 The problem

Objective Bayesian methods, where the formal prior distribution is derived from the assumed model rather than assessed from expert opinions, have a long history (see *e.g.*, Bernardo and Smith, 1994; Kass and Wasserman, 1996, and references therein). Reference priors (Bernardo, 1979, 2005; Berger and Bernardo, 1992a,b, Berger, Bernardo and Sun, 2009, 2012) are a popular choice of objective prior. Other interesting developments involving objective priors include Clarke and Barron (1994), Clarke and Yuan (2004), Consonni, Veronese and Gutiérrez-Peña (2004), DeSantis *et al.* (2001), De Santis (2006), Datta

and Ghosh (1995a, b), Datta and Ghosh (1996), Datta *et al.* (2000), Ghosh (2011), Ghosh, Mergel and Liu (2011), Ghosh and Ramamoorthi (2003), Liseo (1993), Liseo and Loperfido (2006), Sivaganesan (1994), Sivaganesan, Laud and Mueller (2011) and Walker and Gutiérrez-Peña (2011).

In single parameter problems, the reference prior is uniquely defined and is invariant under reparameterization. However, in multiparameter models, the reference prior depends on the quantity of interest, *e.g.*, the parameter concerning which inference is being performed. Thus, if data  $\mathbf{x}$  are assumed to have been generated from  $p(\mathbf{x}|\boldsymbol{\omega})$ , with  $\boldsymbol{\omega} \in \Omega \subset \mathbb{R}^k$ , and one is interested in  $\theta(\boldsymbol{\omega}) \in \Theta \subset \mathbb{R}$ , the reference prior  $\pi_\theta(\boldsymbol{\omega})$ , will typically depend on  $\theta$ ; the posterior distribution,  $\pi_\theta(\boldsymbol{\omega}|\mathbf{x}) \propto p(\mathbf{x}|\boldsymbol{\omega})\pi_\theta(\boldsymbol{\omega})$ , thus also depends on  $\theta$ , and inference for  $\theta$  is performed using the corresponding marginal reference posterior for  $\theta(\boldsymbol{\omega})$ , denoted  $\pi_\theta(\theta|\mathbf{x})$ . The dependence of the reference prior on the quantity of interest has proved necessary to obtain objective posteriors with appropriate properties; in particular, to avoid marginalization paradoxes and strong inconsistencies, and to have good frequentist coverage properties when attainable.

There are however many situations where one is *simultaneously* interested in all the parameters of the model or perhaps in several functions of them. In prediction or decision analysis for instance, all of the parameters of the model may come into play and often none are individually of major interest. Another situation in which having an overall prior would be beneficial is when a user is interested in a non-standard quantity of interest (*e.g.*, a non-standard function of the model parameters), and is not willing or able to formally derive the reference prior for this quantity of interest. Computation can also be a consideration; having to separately do Bayesian computations with a different reference prior for each parameter can be onerous. Finally, when dealing with non-specialists it may be best pedagogically to just present them with one overall objective prior, rather than attempting to explain the technical reasons for preferring different reference priors for different quantities of interest.

To proceed, let  $\boldsymbol{\theta} = \boldsymbol{\theta}(\boldsymbol{\omega}) = \{\theta_1(\boldsymbol{\omega}), \dots, \theta_m(\boldsymbol{\omega})\}$  be the set of  $m > 1$  functions of interest. Our goal is to find a joint prior  $\pi(\boldsymbol{\omega})$  whose corresponding marginal posteriors,  $\{\pi(\theta_i|\mathbf{x})\}_{i=1}^m$ , are sensible from a reference prior perspective. This is not a well-defined goal, and so we will explore various possible approaches to the problem.

**Example 1.1 Multinomial Example:** Suppose  $\mathbf{x} = (x_1, \dots, x_m)$ , where  $\sum_{i=1}^m x_i = n$ , is multinomial  $\text{Mu}(\mathbf{x}|n; \theta_1, \dots, \theta_m)$ , with  $\sum_{i=1}^m \theta_i = 1$ . In Berger and Bernardo (1992b), the reference prior is derived when the parameter  $\theta_i$  is of interest, and this is a different prior for each  $\theta_i$ , as given in the paper. The reference prior for  $\theta_i$  results in a Beta reference marginal posterior  $\text{Be}(\theta_i|x_i + \frac{1}{2}, n - x_i + \frac{1}{2})$ . We would like to identify a single joint prior for  $\boldsymbol{\theta}$  whose marginal posteriors could be expected to be close to each of these reference marginal posteriors, in some average sense.

## 1.2 Background

It is useful to begin by recalling earlier efforts at obtaining an overall reference prior. There have certainly been analyses that can be interpreted as informal efforts at obtaining an overall reference prior. One example is given in Berger and Sun (2008) for the five parameter bivariate normal model. Priors for all the quantities of interest that had previously been considered for the bivariate normal model (21 in all) were studied from a variety of perspectives. One such perspective was that of finding a good overall prior, defined as one which yielded reasonable frequentist coverage properties when used for at least the most important quantities of interest. The conclusion was that the prior  $\pi^o(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho) = 1/[\sigma_1\sigma_2(1 - \rho^2)]$ , where the  $\mu_i$  are the means, the  $\sigma_i$  are the standard deviations, and  $\rho$  is the correlation in the bivariate normal model, was a good choice for the overall prior.

We now turn to some of the more formal efforts to create an overall objective prior.

### 1.2.1 Invariance-based priors

If  $p(\mathbf{x} \mid \boldsymbol{\omega})$  has a group invariance structure, then the recommended objective prior is typically the right-Haar prior. Often this will work well for all parameters that define the invariance structure. For instance, if the sampling model is  $N(x_i \mid \mu, \sigma)$ , the right-Haar prior is  $\pi(\mu, \sigma) = \sigma^{-1}$ , and this is fine for either  $\mu$  or  $\sigma$  (yielding the usual objective posteriors). Such a nice situation does not always obtain, however.

**Example 1.2 Bivariate Normal Distribution:** The right-Haar prior is not unique for the bivariate normal problem. For instance, two possible right-Haar priors are  $\pi_1(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho) = 1/[\sigma_1^2(1 - \rho^2)]$  and  $\pi_2(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho) = 1/[\sigma_2^2(1 - \rho^2)]$ . In Berger and Sun (2008) it is shown that  $\pi_i$  is fine for  $\mu_i$ ,  $\sigma_i$  and  $\rho$ , but leads to problematical posteriors for the other mean and standard deviation.

The situation can be even worse if the right-Haar prior is used for other parameters that can be considered.

**Example 1.3 Multi-Normal Means:** Let  $x_i$  be independent normal with mean  $\mu_i$  and variance 1, for  $i = 1 \dots, m$ . The right-Haar prior for  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)$  is just a constant, which is fine for each of the individual normal means, resulting in a sensible  $N(\mu_i \mid x_i, 1)$  posterior for each individual  $\mu_i$ . But this prior is bad for overall quantities such as  $\theta = \frac{1}{m}|\boldsymbol{\mu}|^2 = \frac{1}{m} \sum_{i=1}^m \mu_i^2$ , as discussed in Stein (1959) and Bernardo and Smith (1994, p. 365). For instance, the resulting posterior mean of  $\theta$  is  $[1 + \frac{1}{m} \sum_{i=1}^m x_i^2]$ , which is inconsistent as  $m \rightarrow \infty$  (assuming  $\frac{1}{m} \sum_{i=1}^m \mu_i^2$  has a limit); indeed, it is easy to show that then  $[1 + \frac{1}{m} \sum_{i=1}^m x_i^2] \rightarrow [\theta_T + 2]$ , where  $\theta_T$  is the true value of  $\theta$ . Furthermore, the posterior distribution of  $\theta$  concentrates sharply around this incorrect value.

### 1.2.2 Constant and vague proper priors

Laplace (1812) advocated use of a constant prior as the overall objective prior and this approach, eventually named *inverse probability*, dominated statistical practice for over 100 years. But the problems of a constant prior are well-documented, including the following:

- (i) Lack of invariance to transformation, the main criticism directed at Laplace's approach.
- (ii) Frequent posterior impropriety.
- (iii) Possible terrible performance, as in the earlier multi-normal mean example.

Vague proper priors (such as a constant prior over a large compact set) are perceived by many as being adequate as an overall objective prior, but they too have well-understood problems. Indeed, they are, at best, equivalent to use of a constant prior, and so inherit most of the flaws of a constant prior. In the multi-normal mean example, for instance, use of  $N(\mu_i \mid 0, 1000)$  vague proper priors results in a posterior mean for  $\theta$  that is virtually identical to the inconsistent posterior mean from the constant prior.

There is a common misperception that vague proper priors are safer than a constant prior, since a proper posterior is guaranteed with a vague proper prior but not for a constant prior. But this actually makes vague proper priors more dangerous than a constant prior. When the constant prior results in an improper posterior distribution, the vague proper prior will yield an essentially arbitrary posterior, depending on the degree of vagueness that is chosen for the prior. And to detect that the answer is arbitrary, one has to conduct a sensitivity study concerning the degree of vagueness, something that can be difficult in complex problems when several or high-dimensional vague proper priors are used. With the constant prior on the other hand, the impropriety of the posterior will usually show up in the computation (the MCMC will not converge) and hence can be recognized.

### 1.2.3 Jeffreys-rule prior

The Jeffreys-rule prior (Jeffreys, 1946, 1961) is the same for all parameters in a model, and is, hence, an obvious candidate for an overall prior. If the data model density is  $p(\mathbf{x} \mid \boldsymbol{\theta})$  the Jeffreys-rule prior for the unknown  $\boldsymbol{\theta} = \{\theta_1, \dots, \theta_m\}$  has the form

$$\pi(\theta_1, \dots, \theta_m) = |I(\boldsymbol{\theta})|^{1/2},$$

where  $I(\boldsymbol{\theta})$  is the  $m \times m$  Fisher information matrix with  $(i, j)$  element

$$I(\boldsymbol{\theta})_{ij} = E_{\mathbf{x} \mid \boldsymbol{\theta}} \left[ - \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log p(\mathbf{x} \mid \boldsymbol{\theta}) \right].$$

This is the optimal objective prior (from many perspectives) for *regular one-parameter* models, but has problems for multi-parameter models. For instance, the right-Haar prior in

the earlier multi-normal mean problem is also the Jeffreys-rule prior there, and was seen to result in an inconsistent estimator of  $\theta$ . Even for the basic  $N(x_i | \mu, \sigma)$  model, the Jeffreys-rule prior is  $\pi(\mu, \sigma) = 1/\sigma^2$ , which results in posterior inferences for  $\mu$  and  $\sigma$  that have the wrong ‘degrees of freedom.’

For the bivariate normal example, the Jeffreys-rule prior is  $1/[\sigma_1^2 \sigma_2^2 (1 - \rho^2)^2]$ ; this yields the natural marginal posteriors for the means and standard deviations, but results in quite inferior objective posteriors for  $\rho$  and various derived parameters, as shown in Berger and Sun (2008). More, generally, the Jeffreys-rule prior for a covariance matrix is studied in Yang and Berger (1994), and shown to yield a decidedly inferior posterior.

There have been efforts to improve upon the Jeffreys-rule prior, such as consideration of the “independence Jeffreys-rule prior,” but such prescriptions have been rather *ad hoc* and have not lead to a general alternative definition.

Finally, consider the following well-known example, which suggests problems with the Jeffreys-rule prior even when it is proper.

**Example 1.4 Multinomial distribution (continued):** Consider the multinomial example where the sample size  $n$  is small relative to the number of classes  $m$ ; thus we have a large sparse table. The Jeffreys-rule prior,  $\pi(\theta_1, \dots, \theta_m) \propto \prod_{i=1}^m \theta_i^{-1/2}$  is a proper prior, but is not a good candidate for the overall prior. For instance, suppose  $n = 3$  and  $m = 1000$ , with  $x_{240} = 2$ ,  $x_{876} = 1$ , and all the other  $x_i = 0$ . The posterior means resulting from use of the Jeffreys-rule prior are

$$E[\theta_i | \mathbf{x}] = \frac{x_i + 1/2}{\sum_{i=1}^m (x_i + 1/2)} = \frac{x_i + 1/2}{n + m/2} = \frac{x_i + 1/2}{503},$$

so  $E[\theta_{240} | \mathbf{x}] = \frac{2.5}{503}$ ,  $E[\theta_{876} | \mathbf{x}] = \frac{1.5}{503}$ ,  $E[\theta_i | \mathbf{x}] = \frac{0.5}{503}$  otherwise. So, cells 240 and 876 only have total posterior probability of  $\frac{4}{503} = 0.008$  even though all 3 observations are in these cells. The problem is that the Jeffreys-rule prior effectively added 1/2 to the 998 zero cells, making them more important than the cells with data! That the Jeffreys-rule prior can encode much more information than is contained in the data is hardly desirable for an objective analysis.

An alternative overall prior that is sometimes considered is the uniform prior on the simplex, but this is even worse than the Jeffreys prior, adding 1 to each cell. The prior that adds 0 to each cell is  $\prod_i \theta_i^{-1}$ , but this results in an improper posterior if any cell has a zero entry, a virtual certainty for very large tables.

We actually know of no multivariable example in which we would recommend the Jeffreys-rule prior. In higher dimensions, the prior always seems to be either ‘too diffuse’ as in the multinormal means example, or ‘too concentrated’ as in the multinomial example.

## 1.3 Three approaches to construction of the overall prior

### 1.3.1 Reference distance approach

In this approach, one seeks a prior that will yield marginal posteriors, for each  $\theta_i$  of interest, that are close to the set of reference posteriors  $\{\pi(\theta_i | \mathbf{x})\}_{i=1}^m$  (yielded by the set of reference priors  $\{\pi_{\theta_i}(\boldsymbol{\omega})\}_{i=1}^m$ ), in an average sense over both posteriors and data  $\mathbf{x} \in \mathcal{X}$ .

**Example 1.5 Multinomial example (continued):** In Example 1.4 consider, as an overall prior, the Dirichlet  $\text{Di}(\boldsymbol{\theta} | a, \dots, a)$  distribution, having density proportional to  $\prod_i \theta_i^{(a-1)}$ . The marginal posterior for  $\theta_i$  is then  $\text{Be}(\theta_i | x_i + a, n - x_i + (m-1)a)$ . In Section 2.2.3, we will study which choice of  $a$  yields marginal posteriors that are as close as possible to the reference marginal posteriors  $\text{Be}(\theta_i | x_i + \frac{1}{2}, n - x_i + \frac{1}{2})$ , arising when  $\theta_i$  is the parameter of interest. Roughly, the recommended choice is  $a = 1/m$ , resulting in the overall prior  $\pi^o(\theta_1, \dots, \theta_m) \propto \prod_{i=1}^m \theta_i^{(1-m)/m}$ . Note that this distribution adds only  $1/m = 0.001$  to each cell in the earlier example, so that

$$\mathbb{E}[\theta_i | \mathbf{x}] = \frac{x_i + 1/m}{\sum_{i=1}^m (x_i + 1/m)} = \frac{x_i + 1/m}{n + 1} = \frac{x_i + 0.001}{4}.$$

Thus  $\mathbb{E}[\theta_{240} | \mathbf{x}] \approx 0.5$ ,  $\mathbb{E}[\theta_{876} | \mathbf{x}] \approx 0.25$ , and  $\mathbb{E}[\theta_i | \mathbf{x}] \approx \frac{1}{4000}$  otherwise, all sensible results.

### 1.3.2 Prior averaging approach

Starting with a collection of reference (or other) priors  $\{\pi_i(\boldsymbol{\theta}), i = 1, \dots, m\}$  for differing parameters or quantities of interest, a rather natural approach is to use an average of the priors. Two natural averages to consider are the arithmetic mean

$$\pi^A(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m \pi_i(\boldsymbol{\theta}),$$

and the geometric mean

$$\pi^G(\boldsymbol{\theta}) = \prod_{i=1}^m \pi_i(\boldsymbol{\theta})^{1/m}.$$

While the arithmetic average might seem most natural, arising from the hierarchical reasoning of assigning each  $\pi_i$  probability  $1/m$  of being correct, geometric averaging arises naturally in the definition of reference priors (Berger, Bernardo and Sun, 2009), and also is the optimal prior if one is trying to choose a single prior to minimize the average of the Kullback-Liebler divergences of the prior from the  $\pi_i$ 's (a fact of which we were reminded by Gauri Datta). Furthermore, the weights in arithmetic averaging of improper priors are rather arbitrary because the priors have no normalizing constants, whereas geometric averaging is unaffected by normalizing constants.

**Example 1.6 Bivariate Normal Distribution (continued):** Faced with the two right-Haar priors in this problem,

$$\pi_1(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho) = \sigma_1^{-2}(1 - \rho^2)^{-1}, \quad \pi_2(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho) = \sigma_2^{-2}(1 - \rho^2)^{-1},$$

the two average priors are

$$\pi^A(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho) = \frac{1}{2\sigma_1^2(1 - \rho^2)} + \frac{1}{2\sigma_2^2(1 - \rho^2)}, \quad (1)$$

$$\pi^G(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho) = \frac{1}{\sigma_1\sigma_2(1 - \rho^2)}. \quad (2)$$

Interestingly, Sun and Berger (2007) show that  $\pi^A$  is a worse objective prior than either right-Haar prior alone, while  $\pi^G$  is the overall recommended objective prior.

One problem with the averaging approach is that the each of the reference priors can depend on all of the other parameters, and not just the parameter of interest,  $\theta_i$ , for which it was created.

**Example 1.7 Multinomial example (continued):** The reference prior derived when the parameter of interest is  $\theta_i$  actually depends on the ordering chosen (in the reference prior derivation) for all the parameters (e.g.  $\{\theta_i, \theta_1, \theta_2, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_m\}$ ); there are thus  $(m - 1)!$  different reference priors for each parameter of interest. Each of these reference priors will result in the same marginal reference posterior for  $\theta_i$ ,

$$\pi_{\theta_i}(\theta_i | \mathbf{x}) = \text{Be}(\theta_i | x_i + \tfrac{1}{2}, n - x_i + \tfrac{1}{2}),$$

but the full reference prior and the full posterior,  $\pi_{\theta_i}(\theta | \mathbf{x})$ , do depend on the ordering of the other parameters. There are thus a total of  $m!$  such full reference priors to be averaged, leading to an often-prohibitive computation.

In general, the quality of reference priors as overall priors is unclear, so there is no obvious sense in which an average of them will make a good overall reference prior. The prior averaging approach is thus best viewed as a method of generating interesting possible priors for further study, and so will not be considered further herein.

### 1.3.3 Hierarchical approach

Utilize hierarchical modeling to transfer the reference prior problem to a ‘higher level’ of the model (following the advice of I. J. Good). In this approach one

- (i) Chooses a class of *proper* priors  $\pi(\boldsymbol{\theta} | a)$  reflecting the desired structure of the problem.
- (ii) Forms the marginal likelihood  $p(\mathbf{x} | a) = \int p(\mathbf{x} | a)\pi(\boldsymbol{\theta} | a) d\boldsymbol{\theta}$ .
- (iii) Finds the reference prior,  $\pi^R(a)$ , for  $a$  in this marginal model.

Thus the overall prior becomes

$$\pi^o(\boldsymbol{\theta}) = \int \pi(\boldsymbol{\theta} | a) \pi^R(a) da,$$

although computation is typically easier in the hierarchical formulation.

**Example 1.8 Multinomial (continued)** The Dirichlet  $\text{Di}(\boldsymbol{\theta} | a, \dots, a)$  class of priors is natural here, reflecting the desire to treat all the  $\theta_i$  similarly. We thus need only to find the reference prior for  $a$  in the marginal model.

$$\begin{aligned} p(\mathbf{x} | a) &= \int \binom{n}{x_1 \dots x_m} \left( \prod_{i=1}^m \theta_i^{x_i} \right) \frac{\Gamma(ma)}{\Gamma(a)^m} \prod_{i=1}^m \theta_i^{a-1} d\boldsymbol{\theta} \\ &= \binom{n}{x_1 \dots x_m} \frac{\Gamma(ma)}{\Gamma(a)^m} \frac{\prod_{i=1}^m \Gamma(x_i + a)}{\Gamma(n + ma)}. \end{aligned} \quad (3)$$

The reference prior for  $\pi^R(a)$  would just be the Jeffreys-rule prior for this marginal model; this is computed in Section 3. The implied prior for  $\boldsymbol{\theta}$  is, of course

$$\pi(\boldsymbol{\theta}) = \int \text{Di}(\boldsymbol{\theta} | a) \pi^R(a) da.$$

Interestingly,  $\pi^R(a)$  turns out to be a proper prior, necessary because the marginal likelihood is bounded away from zero as  $a \rightarrow \infty$ .

As computations in this hierarchical setting are more complex, one might alternatively simply choose the Type-II maximum likelihood estimate, *i.e.*, the value of  $a$  that maximizes (3). For the data given in the earlier example (one cell having two counts, another one count, and the rest zero counts), this marginal likelihood is proportional to  $[a(a+1)]/[(ma+1)(ma+2)]$ , which is maximized at roughly  $a = \sqrt{2}/m$ . In Section 3 we will see that it is actually considerably better to maximize the reference posterior for  $a$ , namely  $\pi^R(a | \mathbf{x}) \propto p(\mathbf{x} | a) \pi^R(a)$ .

## 1.4 Outline of the paper

In Section 2, we formalize the reference distance approach and apply it three models—the multinomial model, the normal model where the coefficient of variation is also a parameter of interest, and a hypergeometric model. In Section 3 we consider the hierarchical prior modeling approach, applying it to three models—the multinomial model, the multinormal model, and the bivariate normal model. Section 4 presents conclusions.



## 2 Reference distance approach

Recall that the goal is to identify a single overall prior  $\pi(\boldsymbol{\omega})$  that can be systematically used for all the parameters  $\boldsymbol{\theta} = \boldsymbol{\theta}(\boldsymbol{\omega}) = \{\theta_1(\boldsymbol{\omega}), \dots, \theta_m(\boldsymbol{\omega})\}$  of interest. The idea of the reference distance approach is to find a  $\pi(\boldsymbol{\omega})$  whose corresponding marginal posteriors,  $\{\pi(\theta_i | \mathbf{x})\}_{i=1}^m$  are close, in an average sense, to the reference posteriors  $\{\pi_i(\theta_i | \mathbf{x})\}_{i=1}^m$  arising from the separate reference priors  $\{\pi_{\theta_i}(\boldsymbol{\omega})\}_{i=1}^m$  derived under the assumption that each of the  $\theta_i$ 's is of interest. In the remainder of the paper,  $\boldsymbol{\theta}$  will equal  $\boldsymbol{\omega}$ , so we will drop  $\boldsymbol{\omega}$  from the notation.

We first consider the situation where the problem has an exact solution.

### 2.1 Exact solution

If one is able to find a single joint prior  $\pi(\boldsymbol{\theta})$  whose corresponding marginal posteriors are precisely equal to the reference posteriors for each of the  $\theta_i$ 's, so that, for all  $\mathbf{x} \in \mathcal{X}$ ,

$$\pi(\theta_i | \mathbf{x}) = \pi_i(\theta_i | \mathbf{x}), \quad i = 1, \dots, m,$$

then it is natural to argue that this should be an appropriate solution to the problem. Notice, however, that there may be many joint priors which satisfy this condition. If the joint reference priors for the  $\theta_i$ 's are all equal, then

$$\pi(\boldsymbol{\theta}) = \pi_{\theta_i}(\boldsymbol{\theta}), \quad i = 1 \dots, m,$$

will obviously satisfy the required condition, and is the overall reference prior.

**Example 2.1** *Univariate normal data.* Consider data  $\mathbf{x}$  which consist of a random sample of normal observations, so that  $p(\mathbf{x} | \boldsymbol{\theta}) = p(\mathbf{x} | \mu, \sigma) = \prod_{i=1}^n N(x | \mu, \sigma)$ , and suppose that one is equally interested in  $\mu$  (or any one-to-one transformation of  $\mu$ ) and  $\sigma$  (or any one-to-one transformation of  $\sigma$ , such as the variance  $\sigma^2$ , or the precision  $\sigma^{-2}$ .) The joint reference prior when any of these is the quantity of interest is known to be the right Haar prior  $\pi_\mu(\mu, \sigma) = \pi_\sigma(\mu, \sigma) = \sigma^{-1}$ , and this is thus an exact solution to the overall prior problem under the reference distance approach.

Interestingly, this prior also works well for making *joint inferences on*  $(\mu, \sigma)$  in that it can be verified that the corresponding joint credible regions for  $(\mu, \sigma)$  have appropriate coverage properties. This does not mean, of course, that the overall prior is necessarily good for any function of the two parameters. For instance, if the quantity of interest is the centrality parameter  $\theta = \mu/\sigma$ , the reference prior is easily found to be  $\pi_\theta(\theta, \sigma) = (1 + \frac{1}{2}\theta^2)^{-1/2}\sigma^{-1}$  (Bernardo, 1979), which is not the earlier overall reference prior.

## 2.2 Reference distance solution

When an exact solution is not possible, it is natural to consider a family of candidate prior distributions,  $\mathcal{F} = \{\pi(\boldsymbol{\theta} | \mathbf{a}), \mathbf{a} \in \mathcal{A}\}$ , and choose, as the overall prior, the distribution from this class which yields marginal posteriors that are closest, in an average sense, to the marginal reference posteriors.

### 2.2.1 Directed logarithmic divergence

It is first necessary to decide how to measure the distance between two distributions. We will actually use a divergence, not a distance, namely the directed logarithmic or Kullback-Leibler divergence (Kullback and Leibler, 1951) given in the following definition.

**Definition 1** *Let  $p(\boldsymbol{\psi})$  be the probability density of a random vector  $\boldsymbol{\psi} \in \Psi$ , and consider an approximation  $p_0(\boldsymbol{\psi})$  with the same or larger support. The directed logarithmic divergence of  $p_0$  from  $p$  is*

$$\kappa\{p_0 | p\} = \int_{\Psi} p(\boldsymbol{\psi}) \log \frac{p(\boldsymbol{\psi})}{p_0(\boldsymbol{\psi})} d\boldsymbol{\psi},$$

*provided that the integral exists.*

The non-negative directed logarithmic divergence  $\kappa\{p_0 | p\}$  is the expected log-density ratio of the true density over its approximation; it is invariant under one-to-one transformations of the random vector  $\boldsymbol{\psi}$ ; and it has an operative interpretation as the amount of information (in natural information units or *nits*) which may be expected to be required to recover  $p$  from  $p_0$ . It was first proposed by Stein (1964) as a loss function and, in a decision-theoretic context, it is often referred to as the *entropy loss*.

### 2.2.2 Weighted logarithmic loss

Suppose the relative importance of the  $\theta_i$  is given by a set of weights  $\{w_1, \dots, w_m\}$ , with  $0 < w_i < 1$  and  $\sum_i w_i = 1$ . A natural default value for these is obviously  $w_i = 1/m$ , but there are many situations where this choice may not be appropriate. To define the proposed criterion, we will also need to utilize the reference prior predictives for  $i = 1, \dots, m$ ,

$$p_{\theta_i}(\mathbf{x}) = \int_{\Theta} p(\mathbf{x} | \boldsymbol{\theta}) \pi_{\theta_i}(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$

**Definition 2** *The best overall prior  $\pi^o(\boldsymbol{\theta})$  within the family  $\mathcal{F} = \{\pi(\boldsymbol{\theta} | \mathbf{a}), \mathbf{a} \in \mathcal{A}\}$  is defined as that which minimizes the weighted average expected logarithmic loss, so that*

$$\begin{aligned}\pi^o(\boldsymbol{\theta}) &= \pi(\boldsymbol{\theta} | \mathbf{a}^*), \quad \mathbf{a}^* = \arg \inf_{\mathbf{a} \in \mathcal{A}} d(\mathbf{a}), \\ d(\mathbf{a}) &= \sum_{i=1}^m w_i \int_{\mathcal{X}} \kappa\{\pi_{\boldsymbol{\theta}_i}(\cdot | \mathbf{x}, \mathbf{a}) | \pi_{\boldsymbol{\theta}_i}(\cdot | \mathbf{x})\} p_{\boldsymbol{\theta}_i}(\mathbf{x}) d\mathbf{x}, \quad \mathbf{a} \in \mathcal{A}.\end{aligned}$$

*This can be rewritten, in terms of the sum of expected risks, as*

$$d(\mathbf{a}) = \sum_{i=1}^m w_i \int_{\boldsymbol{\Theta}} \rho_i(\mathbf{a} | \boldsymbol{\theta}) \pi_{\boldsymbol{\theta}_i}(\boldsymbol{\theta}) d\boldsymbol{\theta}, \quad \mathbf{a} \in \mathcal{A},$$

where

$$\rho_i(\mathbf{a} | \boldsymbol{\theta}) = \int_{\mathcal{X}} \kappa\{\pi_{\boldsymbol{\theta}_i}(\cdot | \mathbf{x}, \mathbf{a}) | \pi_{\boldsymbol{\theta}_i}(\cdot | \mathbf{x})\} p(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x}, \quad \boldsymbol{\theta} \in \boldsymbol{\Theta}.$$

Note that there is no assurance that  $d(\mathbf{a})$  will be finite if the reference priors are improper. Indeed, in cases we have investigated with improper reference priors,  $d(\mathbf{a})$  has failed to be finite and hence the reference distance approach cannot be directly used. However, as in the construction of reference priors, one can consider an approximating sequence of proper priors  $\{\pi_{\boldsymbol{\theta}_i}(\boldsymbol{\theta} | k), k = 1, 2, \dots\}$  on increasing compact sets. For each of the  $\pi_{\boldsymbol{\theta}_i}(\boldsymbol{\theta} | k)$ , one can minimize the expected loss

$$d(\mathbf{a} | k) = \sum_{i=1}^m w_i \int_{\boldsymbol{\Theta}} \rho_i(\mathbf{a} | \boldsymbol{\theta}) \pi_{\boldsymbol{\theta}_i}(\boldsymbol{\theta} | k) d\boldsymbol{\theta},$$

obtaining  $\mathbf{a}_k^* = \arg \inf_{\mathbf{a} \in \mathcal{A}} d(\mathbf{a} | k)$ . Then, if  $\mathbf{a}^* = \lim_{k \rightarrow \infty} \mathbf{a}_k^*$  exists, one can declare this to be the solution.

### 2.2.3 Multinomial model

In the multinomial model with  $m$  cells and parameters  $\{\theta_1, \dots, \theta_m\}$ , with  $\sum_{i=1}^m \theta_i = 1$ , the reference posterior for each of the  $\theta_i$ 's is  $\pi_i(\theta_i | \mathbf{x}) = \text{Be}(\theta_i | x_i + \frac{1}{2}, n - x_i + \frac{1}{2})$ , while the marginal posterior distribution of  $\theta_i$  resulting from the joint prior  $\text{Di}(\theta_1, \dots, \theta_{m-1} | a)$  is  $\text{Be}(\theta_i | x_i + a, n - x_i + (m-1)a)$ . The directed logarithmic discrepancy of the posterior  $\text{Be}(\theta_i | x_i + a, n - x_i + (m-1)a)$  from the reference posterior  $\text{Be}(\theta_i | x_i + \frac{1}{2}, n - x_i + \frac{1}{2})$  is

$$\kappa_i\{a | \mathbf{x}, m, n\} = \kappa_i\{a | x_i, m, n\} = \kappa_{\text{Be}}\{x_i + a, n - x_i + (m-1)a | x_i + \frac{1}{2}, n - x_i + \frac{1}{2}\}$$

where

$$\begin{aligned}
\kappa_{\text{Be}}\{\alpha_0, \beta_0 \mid \alpha, \beta\} &= \int_0^1 \text{Be}(\theta_i \mid \alpha, \beta) \log \left[ \frac{\text{Be}(\theta_i \mid \alpha, \beta)}{\text{Be}(\theta_i \mid \alpha_0, \beta_0)} \right] d\theta_i \\
&= \log \left[ \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha_0 + \beta_0)} \frac{\Gamma(\alpha_0)}{\Gamma(\alpha)} \frac{\Gamma(\beta_0)}{\Gamma(\beta)} \right] \\
&\quad + (\alpha - \alpha_0)\psi(\alpha) + (\beta - \beta_0)\psi(\beta) - ((\alpha + \beta) - (\alpha_0 + \beta_0))\psi(\alpha + \beta),
\end{aligned}$$

and  $\psi(\cdot)$  is the digamma function.

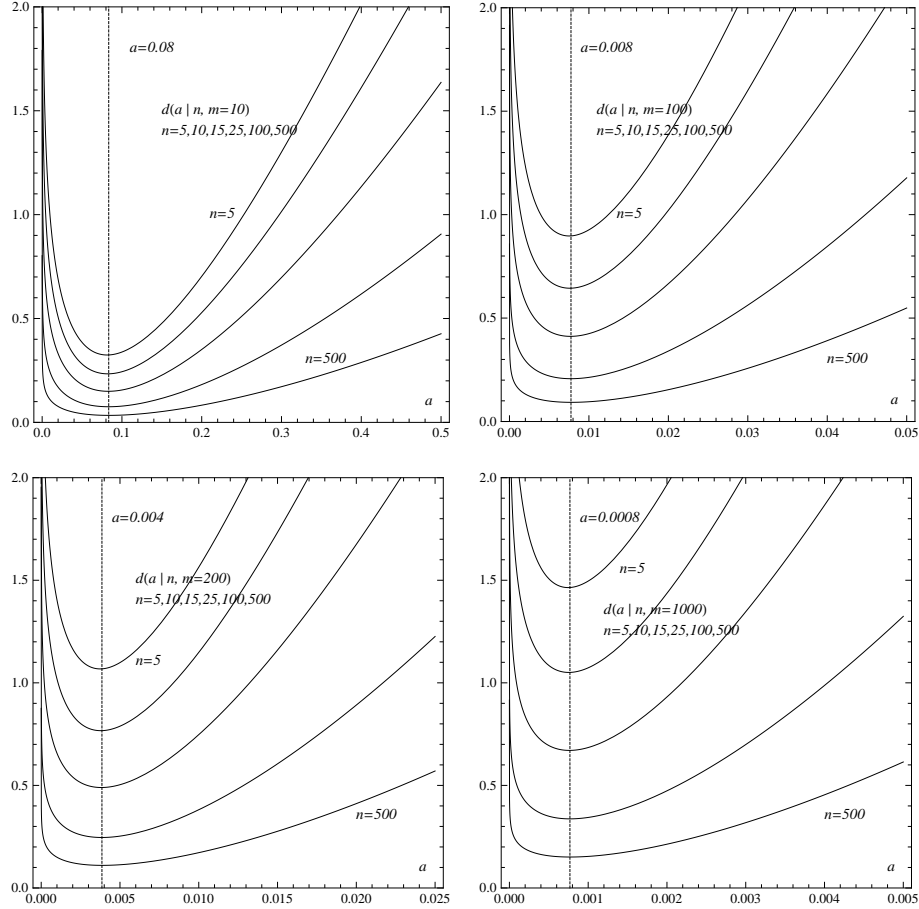


Figure 1: *Expected logarithmic losses, when using a Dirichlet prior with parameter  $\{a, \dots, a\}$ , in a multinomial model with  $m$  cells, for sample sizes  $n = 5, 10, 25, 100$  and 500. Clockwise panels from upper left,  $m = 10, 100, 200$  and 1000. In all cases, the optimal value for all sample sizes is  $a^* \approx 0.8/m$ .*

The divergence  $\kappa_i\{a \mid x_i, m, n\}$  between the two posteriors of  $\theta_i$  depends on the data only through  $x_i$  and the sampling distribution of  $x_i$  is  $\text{Bi}(x_i \mid n, \theta_i)$ , which only depends of  $\theta_i$ . Moreover, the marginal reference prior for  $\theta_i$  is  $\pi_{\theta_i}(\theta_i) = \text{Be}(\theta_i \mid 1/2, 1/2)$  and, therefore, the

corresponding reference predictive for  $x_i$  is

$$p(x_i | n) = \int_0^1 \text{Bi}(x_i | n, \theta_i) \text{Be}(\theta_i | 1/2, 1/2) d\theta_i = \frac{1}{\pi} \frac{\text{Ga}(x_i + \frac{1}{2}) \text{Ga}(n - x_i + \frac{1}{2})}{\text{Ga}(x_i + 1) \text{Ga}(n - x_i + 1)}.$$

Hence, using Definition 2 with uniform weights, the average expected logarithmic loss of using a joint Dirichlet prior with parameter  $a$  with a sample of size  $n$  is simply

$$d(a | m, n) = \sum_{x=0}^n \kappa\{a | x, m, n\} p(x | n)$$

since, by the symmetry of the problem, the  $m$  parameters  $\{\theta_1, \dots, \theta_m\}$  yield the same expected loss.

The function  $d(a | m = 10, n)$  is graphed in the upper left panel of Figure 1 for several values of  $n$ . The expected loss decreases with  $n$  and, for any  $n$ , the function  $d(a | m, n)$  is concave, with a unique minimum numerically found to be at  $a^* \approx 0.8/m = 0.08$ . The approximation is rather precise. For instance, the minimum is achieved at 0.083 for  $n = 100$ .

Similarly, the function  $d(a | m = 1000, n)$  is graphed in the lower right panel of Figure 1 for the same values of  $n$  and with the same vertical scale, yielding qualitatively similar results although, as one may expect, the expected losses are now larger than those obtained with  $m = 10$ . Once more, the function  $d(a | m = 1000, n)$  is concave, with a unique minimum numerically found to be at  $a^* \approx 0.8/m = 0.0008$ , with the exact value very close. For instance, for  $n = 100$ , the minimum is achieved at 0.00076.

It can be concluded that, for all practical purposes when using the reference distance approach, the best global Dirichlet prior, when one is interested in all the parameters of a multinomial model, is that with parameter vector  $\{1/m, \dots, 1/m\}$  (or  $0.8 \times \{1/m, \dots, 1/m\}$  to be slightly more precise), yielding an approximate marginal reference posterior for each of the  $\theta_i$ 's as  $\text{Be}(\theta_i | x_i + 1/m, n - x_i + (m - 1)/m)$ , having mean and variance

$$\mathbb{E}[\theta_i | x_i, n] = \hat{\theta}_i = (x_i + 1/m)/(n + 1), \quad \text{Var}[\theta_i | x_i, n] = \hat{\theta}_i(1 - \hat{\theta}_i)/(n + 2).$$

#### 2.2.4 Multivariate hypergeometric model

Let  $\mathcal{N}_+$  be the set of all nonnegative integers. Consider a multivariate hypergeometric distribution  $\text{Hy}_k(\mathbf{r}_k | n, \mathbf{R}, N)$  with the probability mass function

$$\text{Hy}_k(\mathbf{r}_k | n, \mathbf{R}_k, N) = \frac{\binom{R_1}{r_1} \dots \binom{R_k}{r_k} \binom{R_{k+1}}{r_{k+1}}}{\binom{N}{n}}, \quad \mathbf{r}_k \in \mathcal{R}_{k,n}, \quad (4)$$

$$\mathcal{R}_{k,n} = \{\mathbf{r}_k = (r_1, \dots, r_k); \quad r_j \in \mathcal{N}_+, \quad r_1 + \dots + r_k \leq n\},$$

where the  $k$  unknown parameters  $\mathbf{R}_k = (R_1, \dots, R_k)$  are in the parameter space  $\mathcal{R}_{k,N}$ . Here and in the following,  $R_{k+1} = N - (R_1 + \dots + R_k)$ . Notice that the univariate hypergeometric distribution is the special case when  $k = 1$ .

A natural hierarchical model for the unknown  $\mathbf{R}_k$  is to assume that it is multinomial  $\text{Mu}_k(\mathbf{R}_k | N, \mathbf{p}_k)$ , with  $\mathbf{p}_k \in \mathcal{P}_k \equiv \{\mathbf{p}_k = (p_1, \dots, p_k)\}$ ,  $0 \leq p_j \leq 1$ , and  $p_1 + \dots + p_k \leq 1$ . The probability mass function of  $\mathbf{R}_k$  is then

$$\text{Mu}_k(\mathbf{R}_k | N, \mathbf{p}_k) = \frac{N!}{\prod_{j=1}^{k+1} R_j!} \prod_{j=1}^{k+1} p_j^{R_j}.$$

Berger, Bernardo and Sun (2012) prove that the marginal likelihood of  $\mathbf{r}_k$  given  $(\mathbf{p}_k, n, N)$  depends only on  $(n, \mathbf{p}_k)$  and it is given by

$$\begin{aligned} p(\mathbf{r}_k | \mathbf{p}_k, n, N) &= \sum_{\mathbf{R}_k \in \mathcal{N}_{k,N}} \text{Hy}_k(\mathbf{r}_k | n, \mathbf{R}_k, N) \text{Mu}_k(\mathbf{R}_k | N, \mathbf{p}_k) \\ &= \text{Mu}_k(\mathbf{r}_k | n, \mathbf{p}_k), \quad \mathbf{r}_k \in \mathcal{R}_{k,n}. \end{aligned} \quad (5)$$

This reduces to the multinomial problem. Hence, the overall (approximate) reference prior for  $(\mathbf{R}_k | N, \mathbf{p}_k)$  would be Dirichlet  $\text{Di}(\mathbf{R}_k | 1/k, \dots, 1/k)$ .

### 2.2.5 The normal model with coefficient of variation

Consider a random sample  $\mathbf{z} = \{x_1, \dots, x_n\}$  from a normal model  $N(x | \mu, \sigma)$ , with both parameters unknown, and suppose that one is interested in  $\mu$  and  $\sigma$ , but also in the standardized mean  $\phi = \mu/\sigma$  (and/or any one-to-one function of them such as  $\log \sigma$ , or the coefficient of variation  $\sigma/\mu$ ).

The joint reference prior when either  $\mu$  or  $\sigma$  are the quantities of interest is

$$\pi_\mu(\mu, \sigma) = \pi_\sigma(\mu, \sigma) = \sigma^{-1} \quad (6)$$

and this is known to lead to the reference posteriors

$$\pi_\mu^{ref}(\mu | \mathbf{z}) = \text{St}(\mu | \bar{x}, s/\sqrt{n-1}, n-1), \quad \pi_\sigma^{ref}(\sigma | \mathbf{z}) = \text{Ga}^{-1/2}(\sigma | (n-1)/2, ns^2/2),$$

with  $n\bar{x} = \sum_{i=1}^n x_i$  and  $ns^2 = \sum_{i=1}^n (x_i - \bar{x})^2$ , which are proper if  $n \geq 2$ , and have the correct probability matching properties. However, the reference prior if  $\phi$  is the parameter of interest is  $\pi_\phi(\phi, \sigma) = (2 + \phi^2)^{-1/2} \sigma^{-1}$  (Bernardo, 1979), and the corresponding reference posterior distribution for  $\phi$  can be shown to be

$$\pi_\phi^{ref}(\phi | \mathbf{z}) = \pi_\phi^{ref}(\phi | t) \propto (2 + \phi^2)^{-1/2} p(t | \phi),$$

where  $t = (\sum_{i=1}^n x_i)/(\sum_{i=1}^n x_i^2)^{1/2}$  has a sampling distribution  $p(t | \phi)$  depending only on  $\phi$  (see Stone and Dawid, 1972). Note that all posteriors can be written in terms of the sufficient statistics  $\bar{x}$  and  $s^2$  and the sample size  $n$ , which we will henceforth use.

A natural choice for the family of joint priors to be considered as candidates for an overall prior is the class of *relatively invariant* priors (Hartigan, 1964),

$$\mathcal{F} = \{\pi(\mu, \sigma | a) = \sigma^{-a}, a > 0\}$$

which contains, for  $a = 1$ , the joint reference prior (6) when either  $\mu$  or  $\sigma$  are the parameters of interest, and the Jeffreys-rule prior, for  $a = 2$ . Since these priors are improper, a compact approximation procedure, as described at the end of Section 2.2.2, is needed. The usual compactification for location-scale parameters considers the sets

$$\mathcal{C}_k = \{\mu \in (-k, k), \sigma \in (e^{-k}, e^k)\}, \quad k = 1, 2, \dots$$

One must therefore derive

$$d(a | n, k) = d_\mu(a | n, k) + d_\sigma(a | n, k) + d_\phi(a | n, k),$$

where each of the  $d_i$ 's is found by integrating the corresponding risk with the appropriately renormalized joint reference prior. Thus,

$$\begin{aligned} d_\mu(a | n, k) &= \int_{\mathcal{C}_k} \left[ \int_{\mathcal{T}} \kappa\{\pi_\mu(\cdot | n, \mathbf{t}, a) | \pi_\mu^{ref}(\cdot | n, \mathbf{t})\} p(\mathbf{t} | n, \mu, \sigma) d\mathbf{t} \right] \pi_\mu(\mu, \sigma | k) d\mu d\sigma, \\ d_\sigma(a | n, k) &= \int_{\mathcal{C}_k} \left[ \int_{\mathcal{T}} \kappa\{\pi_\sigma(\cdot | n, \mathbf{t}, a) | \pi_\sigma^{ref}(\cdot | n, \mathbf{t})\} p(\mathbf{t} | n, \mu, \sigma) d\mathbf{t} \right] \pi_\sigma(\mu, \sigma | k) d\mu d\sigma, \\ d_\phi(a | n, k) &= \int_{\mathcal{C}_k} \left[ \int_{\mathcal{T}} \kappa\{\pi_\phi(\cdot | n, \mathbf{t}, a) | \pi_\phi^{ref}(\cdot | n, \mathbf{t})\} p(\mathbf{t} | n, \mu, \sigma) d\mathbf{t} \right] \pi_\phi(\mu, \sigma | k) d\mu d\sigma, \end{aligned}$$

where  $\mathbf{t} = (\bar{x}, s)$ , and the  $\pi_i(\mu, \sigma | k)$ 's are the joint *proper* prior reference densities of each of the parameter functions obtained by truncation and renormalization in the  $\mathcal{C}_k$ 's.

It is found that the risk associated to  $\mu$  (the expected KL divergence of  $\pi_\mu(\cdot | n, \mathbf{t}, a)$  from  $\pi_\mu^{ref}(\cdot | n, \mathbf{t})$  under sampling) does *not* depend on the parameters, so integration with the joint prior is not required, and one obtains

$$d_\mu(a | n, k) = d_\mu(a | n) = \log \left[ \frac{\Gamma[n/2] \Gamma[(a+n)/2 - 1]}{\Gamma[(n-1)/2] \Gamma[(a+n-1)/2]} \right] - \frac{a-1}{2} \left( \psi\left[\frac{n-1}{2}\right] - \psi\left[\frac{n}{2}\right] \right),$$

where  $\psi[\cdot]$  is the digamma function. This is a concave function with a unique minimum  $d_1(1 | n) = 0$  at  $a = 1$ , as one would expect from the fact that the target family  $\mathcal{F}$  contains the reference prior for  $\mu$  when  $a = 1$ . The function  $d_\mu(a | n = 10)$  is the lower dotted line in

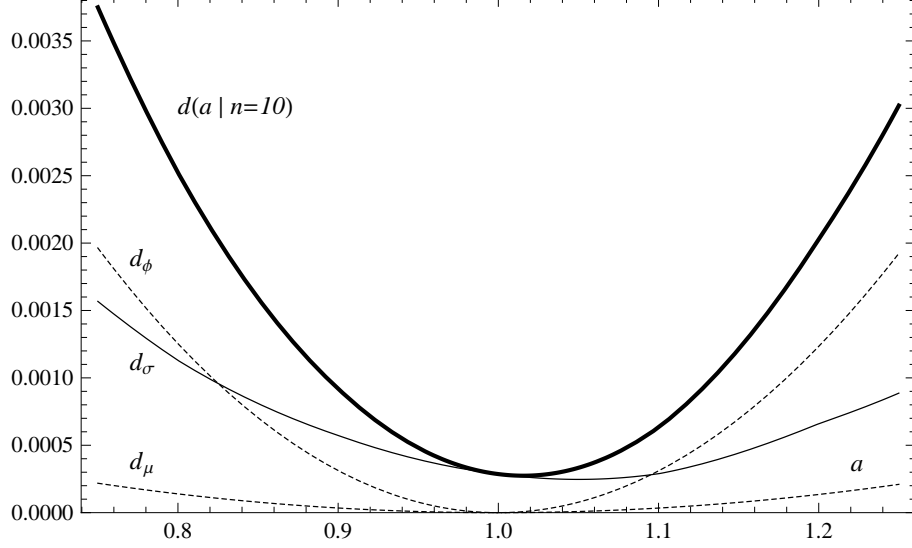


Figure 2: *Expected average intrinsic logarithmic losses  $d(a | n, k)$  associated with the use of the joint prior  $\pi(\mu, \sigma | a) = \sigma^{-a}$  rather than the corresponding reference priors when  $n = 10$  and  $k = 3$ .*

Figure 2. Similarly, the risk associated to  $\sigma$  does not depend either of the parameters, and one obtains

$$d_\sigma(a | n, k) = d_\sigma(a | n) = \log \left[ \frac{\Gamma[(a + n)/2 - 1]}{\Gamma[(n - 1)/2]} \right] - \frac{a - 1}{2} \psi \left[ \frac{n - 1}{2} \right],$$

another concave function which a unique minimum  $d_2(1 | n) = 0$ , at  $a = 1$ . The function  $d_\sigma(a | n = 10)$  is the upper dotted line in Figure 2.

The risk associated with  $\phi$  cannot be analytically obtained and is numerically computed, using one-dimensional numerical integration over  $\phi$  to obtain the KL divergence, and Monte Carlo sampling to obtain its expected value with the truncated and renormalized reference prior  $\pi_\phi(\mu, \sigma | k)$ . The function  $d_\phi(a | n = 10, k = 3)$  is represented by the black line in Figure 2. It may be appreciated that, of the three components of the expected loss, the contribution corresponding to  $\sigma$  is the largest, and that corresponding to  $\mu$  is the smallest, in the neighborhood of the optimal choice of  $a$ . The sum of the three is the expected loss to be minimized,  $d(a | n, k)$ . The function  $d(a | n = 10, k = 3)$  is represented by the solid line in Figure 2, and has a minimum at  $a_3^* = 1.016$ . The sequence of numerically computed optimum values is  $\{a_k^*\} = \{1.139, 1.077, 1.016, \dots\}$  quickly converging to some value  $a^*$  larger than 1 and smaller than 1.016, so that, pragmatically, the overall objective prior may be taken to be the usual objective prior for the normal model,

$$\pi^o(\mu, \sigma) = \sigma^{-1}.$$

It is of interest to study the difference in use of this overall prior when compared with



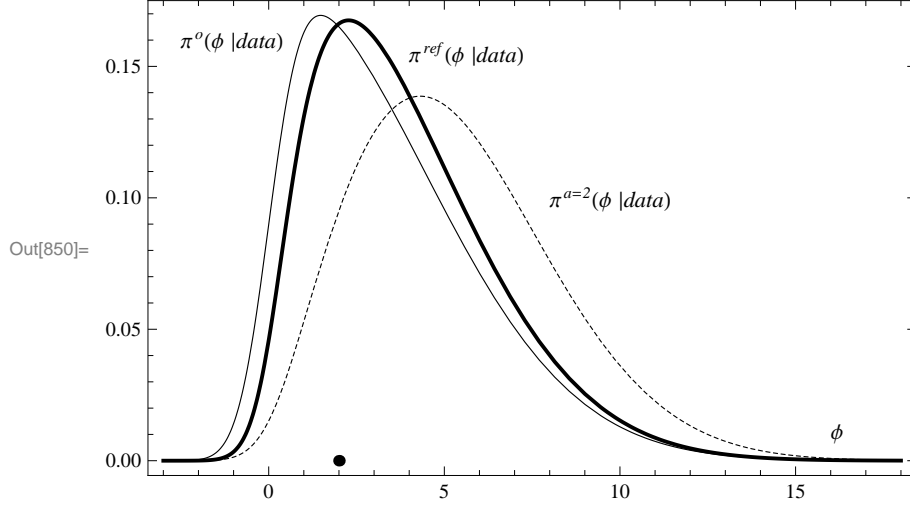


Figure 3: *Reference posterior (solid) and marginal overall posterior (black) for  $\phi$  given a minimal random sample of size  $n = 2$ . The dotted line is the marginal posterior for the prior with  $a = 2$ , which is the Jeffreys-rule prior.*

the reference prior for  $\phi = \mu/\sigma$ . The difference is greater for smaller samples, and the minimum sample size here is  $n = 2$ . A random sample of two observations from  $N(x | 1, \frac{1}{2})$  (so that the true value of the standardized mean is  $\phi = 2$ ) was simulated yielding  $\{x_1, x_2\} = \{0.959, 1.341\}$ . The corresponding reference posterior for  $\phi$  is the solid line in Figure 3. The posterior that corresponds to the recommended overall prior  $a = 1$  is the black line in the figure. For comparison, the posterior corresponding to the prior with  $a = 2$ , which is Jeffreys-rule prior, is also given, as the dotted line. Thus, even with a minimum sample size, the overall prior yields a marginal posterior for  $\phi$  which is quite close to that for the reference posterior. For sample sizes beyond  $n = 4$  the differences are visually inappreciable.

### 3 Hierarchical approach with hyperpriors

If a natural family of proper priors  $\pi(\boldsymbol{\theta} | a)$ , indexed by a single parameter  $a$ , can be identified for a given problem, one can compute the marginal likelihood  $p(\boldsymbol{x} | a)$  (necessarily a proper density), and find the reference prior  $\pi^R(a)$  for  $a$  for this marginal likelihood. This hierarchical prior specification is clearly equivalent to use of

$$\pi^o(\boldsymbol{\theta}) = \int \pi(\boldsymbol{\theta} | a) \pi^R(a) da$$

as the overall prior in the original problem.

## 3.1 Multinomial problem

### 3.1.1 The hierarchical prior

For the multinomial problem with the  $\text{Di}(\boldsymbol{\theta} | a, \dots, a)$  prior, the marginal density of any of the  $x_i$ 's is

$$p(x_i | a, m, n) = \binom{n}{x_i} \frac{\Gamma(x_i + a) \Gamma(n - x_i + (m - 1)a) \Gamma(ma)}{\Gamma(a) \Gamma((m - 1)a) \Gamma(n + ma)},$$

following immediately from the fact that, marginally,

$$p(x_i | \theta_i) = \text{Bi}(x_i | n, \theta_i) \quad \pi(\theta_i | a) = \text{Be}(\theta_i | a, (m - 1)a).$$

Then  $\pi^R(a)$ , the reference (Jeffreys) prior for the integrated model  $p(\mathbf{x} | a)$  in (3), is given in the following proposition:

**Proposition 3.1**

$$\pi^R(a | m, n) \propto \left[ \sum_{j=0}^{n-1} \left( \frac{Q(j | a, m, n)}{(a + j)^2} - \frac{m}{(ma + j)^2} \right) \right]^{1/2}, \quad (7)$$

where  $Q(\cdot | a, m, n)$  is the right tail of the distribution of  $p(x | a, m, n)$ , namely

$$Q(j | a, m, n) = \sum_{l=j+1}^n p(l | a, m, n), \quad j = 0, \dots, n - 1.$$

*Proof.* Computation yields that

$$\mathbb{E} \left[ -\frac{d^2}{da^2} \log p(\mathbf{x} | a) \right] = -\sum_{j=0}^{n-1} \frac{m^2}{(ma + j)^2} + E \left[ \sum_{i=1}^m \sum_{j=0}^{x_i-1} \frac{1}{(a + j)^2} \right], \quad (8)$$

where  $\sum_{j=0}^{-1} \equiv 0$ . Since the  $x_i$  are exchangeable, this equals

$$-\sum_{j=0}^{n-1} \frac{m^2}{(ma + j)^2} + mE^{X_1} \left[ \sum_{j=0}^{X_1-1} \frac{1}{(a + j)^2} \right],$$

and the result follows by rearranging terms.  $\square$

**Proposition 3.2**  $\pi^R(a)$  is a proper prior.

*Proof.* The prior is clearly continuous in  $a$ , so we only need show that it is integrable at 0 and at  $\infty$ . Consider first the situation as  $a \rightarrow \infty$ . Then

$$\begin{aligned} p(0 | a, m, n) &= \frac{\Gamma(a) \Gamma(n + [m - 1]a) \Gamma(ma)}{\Gamma(a) \Gamma([m - 1]a) \Gamma(n + ma)} \\ &= \frac{(m - 1)a[(m - 1)a + 1] \cdots [(m - 1)a + n - 1]}{ma(ma + 1) \cdots (ma + n - 1)} \\ &= \frac{(m - 1)}{m} (1 - c_n a + O(a^2)), \end{aligned}$$

where  $c_n = 1 + 1/2 + \dots + 1/(n-1)$ . Thus the first term of the sum in (7) is

$$\frac{1 - p(0 | a, m, n)}{a^2} - \frac{1}{m a^2} = \frac{(m-1)c_n}{m a} + O(1).$$

All of the other terms of the sum in (7) are clearly  $O(1)$ , so that

$$\pi^R(a) = \frac{\sqrt{(m-1)c_n/m}}{\sqrt{a}} + O(\sqrt{a}),$$

as  $a \rightarrow 0$ , which is integrable at zero (although unbounded).

To study propriety as  $a \rightarrow \infty$ , a laborious application of Stirling's approximation yields

$$p(x_1 | a, m, n) = \text{Bi}(x_1 | n, 1/m)(1 + O(a^{-1})),$$

as  $a \rightarrow \infty$ . Thus

$$\begin{aligned} \pi^R(a, m, n) &= \left[ \sum_{j=0}^{n-1} \left( \frac{\sum_{l=j+1}^n \text{Bi}(l | n, 1/m)}{a^2} - \frac{1}{m a^2} \right) + O(a^{-3}) \right]^{1/2} \\ &= \left[ \left( \frac{\sum_{l=1}^n l \text{Bi}(l | n, 1/m)}{a^2} - \frac{n}{m a^2} \right) + O(a^{-3}) \right]^{1/2} = O(a^{-3/2}), \end{aligned}$$

which is integrable at infinity, completing the proof.  $\square$

As suggested by the proof above, the reference prior  $\pi^R(a | m, n)$  behaves as  $O(a^{-1/2})$  near  $a = 0$  and behaves as  $O(a^{-2})$  for large  $a$  values. Using series expansions, it is found that, for sparse tables where  $m/n$  is relatively large, the reference prior is well approximated by the proper prior

$$\pi^*(a | m, n) = \frac{1}{2} \frac{n}{m} a^{-1/2} \left( a + \frac{n}{m} \right)^{-3/2}, \quad (9)$$

which only depends on the ratio  $m/n$ , and has the behavior at the extremes described above. This can be restated as saying that  $\phi(a) = a/(a + (n/m))$  has a Beta distribution  $\text{Be}(\phi | \frac{1}{2}, 1)$ . Figure 4 gives the exact form of  $\pi^R(a | m, n)$  for various  $(m, n)$  values, and the corresponding approximation given by (9). The approximate reference prior  $\pi^*(a | m, n)$  appears to be a good approximation to the actual reference prior, and hence can be recommended for use with large sparse contingency tables.

It is always a surprise when a reference prior turns out to be proper, and this seems to happen when the likelihood does not go to zero at a limit. Indeed, it is straightforward to show that

$$p(\mathbf{x} | a) = \begin{cases} O(a^{r_0-1}), & \text{as } a \rightarrow 0, \\ \binom{n}{\mathbf{x}} m^{-n}, & \text{as } a \rightarrow \infty, \end{cases}$$

where  $r_0$  is the number of nonzero  $x_i$ . Thus, indeed, the likelihood is constant at  $\infty$ , so that the prior must be proper at infinity for the posterior to exist.

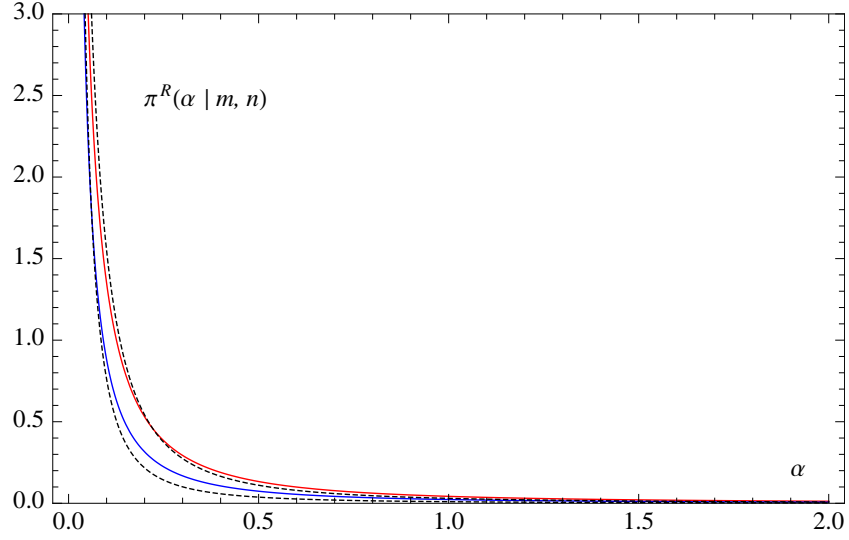


Figure 4: Reference priors  $\pi^R(a | m, n)$  (solid lines) and their approximations (dotted lines) for  $(m = 150, n = 10)$  (upper curve) and for  $(m = 500, n = 10)$  (lower curve).

### 3.1.2 Computation with the hierarchical reference prior

If a full Bayesian analysis is desired, the obvious MCMC sampler is as follows:

*Step 1.* Use a Metropolis Hastings move to sample from the marginal posterior  $\pi^R(a | \mathbf{x}) \propto \pi^R(a) p(\mathbf{x} | a)$ .

*Step 2.* Given  $a$ , sample from the usual beta posterior  $\pi(\theta | a, \mathbf{x})$ .

This will be highly efficient if a good proposal distribution for Step 1 can be found. As it is only a one-dimensional distribution, standard techniques should work well. Even simpler computationally is the use of the approximate reference prior  $\pi^*(a | m, n)$  in (9), because of the following result.

**Proposition 3.3** *Under the approximate reference prior (9), and provided there are at least three nonempty cells, the marginal posterior distribution of  $a$  is log-concave.*

*Proof.* It follows from (8) that

$$\frac{d^2}{da^2} \log[p(\mathbf{x} | a) \pi^*(a | m, n)] = \sum_{j=0}^{n-1} \frac{m^2}{(ma + j)^2} - \sum_{i=1}^m \sum_{j=0}^{x_i-1} \frac{1}{(a + j)^2} + \frac{1}{2a^2} + \frac{3}{2(a + n/m)^2}.$$

Without loss of generality, we assume that  $x_i > 0$ , for  $i = 1, 2, 3$ . Then

$$\frac{d^2}{da^2} \log[p(\mathbf{x} | a) \pi^*(a | m, n)] < - \sum_{i=2}^3 \sum_{j=0}^{x_i-1} \frac{1}{(a + j)^2} + \frac{1}{2a^2} + \frac{3}{2a^2} < 0. \quad \square$$

Thus log-concave rejection sampling (Gilks and Wild, 1992) can be used to sample from the posterior of  $a$ .

Alternatively, one might consider the empirical Bayes solution of fixing  $a$  at its posterior mode  $\hat{a}^R$ . The one caveat is that, when  $r_0 = 1$ , it follows from (10) that the likelihood is constant at zero, while  $\pi^R(a)$  is unbounded at zero; hence the posterior mode will be  $a = 0$ , which cannot be used. When  $r_0 \geq 2$ , it is easy to see that  $\pi^R(a)p(\mathbf{x}|a)$  goes to zero as  $a \rightarrow 0$ , so there will be no problem.

It will typically be considerably better to utilize the posterior mode than the maximum of  $p(\mathbf{x}|a)$  alone, given the fact that the likelihood does not go to zero at  $\infty$ . For instance, if all  $x_i = 1$ , it can be shown that  $p(\mathbf{x}|a)$  has a likelihood increasing in  $a$ , so that there is no mode. (Even when  $r_0 = 1$ , use of the mode of  $p(\mathbf{x}|a)$  is not superior, in that the likelihood is also maximized at 0 in that case.)

### 3.1.3 Posterior behavior as $m \rightarrow \infty$

Since we are contemplating the “large sparse” contingency table scenario, it is of considerable interest to study the behavior of the posterior distribution as  $m \rightarrow \infty$ . It is easiest to state the result in terms of the transformed variable  $v = ma$ . Let  $\pi_m^R(v|\mathbf{x})$  denote the transformed reference posterior.

#### Proposition 3.4

$$\Psi(v) = \lim_{m \rightarrow \infty} \pi_m^R(v|\mathbf{x}) = \frac{\Gamma(v+1)}{\Gamma(v+n)} v^{(r_0 - \frac{3}{2})} \left[ \sum_{i=1}^{n-1} \frac{i}{(v+i)^2} \right]^{1/2}. \quad (10)$$

*Proof.* Note that

$$\begin{aligned} \pi^R(a|\mathbf{x}) &\propto m(\mathbf{x}|a) \pi^R(a) \\ &\propto \frac{\Gamma(ma)}{\Gamma(ma+n)} \left[ \prod_{i=1}^m \frac{\Gamma(a+x_i)}{\Gamma(a)} \right] \pi^R(a) \\ &\propto \frac{\Gamma(ma)}{\Gamma(ma+n)} \left[ \prod_{i:x_i \neq 0} a(a+1) s(a+x_i-1) \right] \pi^R(a) \\ &\propto \frac{\Gamma(ma)}{\Gamma(ma+n)} \left[ \prod_{j=0}^{n-1} (a+j)^{r_j} \right] \pi^R(a), \end{aligned}$$

where  $r_j = \{\#x_i > j\}$ . Change of variables to  $v = ma$  yields

$$\begin{aligned} \pi_m^R(v|\mathbf{x}) &\propto \frac{\Gamma(v)}{\Gamma(v+n)} \left[ \prod_{j=0}^{n-1} \left( \frac{v}{m} + j \right)^{r_j} \right] \pi^R\left(\frac{v}{m}\right) \\ &\propto \frac{\Gamma(v)}{\Gamma(v+n)} v^{r_0} \left[ C + \sum_{i=1}^{n-r_0} K_i \left( \frac{v}{m} \right)^i \right] \pi^R\left(\frac{v}{m}\right), \end{aligned} \quad (11)$$

where  $C = \prod_{j=2}^{n-1} j^{r_j}$  and the  $K_i$  are constants.

Next we study the behavior of  $\pi^R(v/m)$  for large  $m$ . Note first that, in terms of  $v$ , the marginal density of  $x_1 = 0$  is

$$\begin{aligned} p(0|v) &= \frac{\Gamma(\frac{(m-1)}{m}v + n)}{\Gamma(\frac{(m-1)}{m}v)} \frac{\Gamma(v)}{\Gamma(v+n)} \\ &= \frac{\frac{(m-1)}{m}v [\frac{(m-1)}{m}v + 1] \cdots [\frac{(m-1)}{m}v + n - 1]}{v(v+1) \cdots (v+n-1)} \\ &= \frac{(m-1)}{m} \left(1 - \frac{v}{m[v+1]}\right) \cdots \left(1 - \frac{v}{m[v+n-1]}\right) \\ &= \frac{(m-1)}{m} \left(1 - \frac{v}{m} \sum_{i=1}^{n-1} \frac{1}{v+i} + O\left(\frac{v^2}{m^2(v+1)^2}\right)\right). \end{aligned}$$

Hence

$$Q(0|a) = 1 - p(0|v) = \frac{1}{m} + \frac{v(m-1)}{m^2} \sum_{i=1}^{n-1} \frac{1}{v+i} + O\left(\frac{v^2}{m^2(v+1)^2}\right) = O\left(\frac{1}{m}\right) \text{ (uniformly in } v\text{)}.$$

It follows that all  $Q(i|a) \leq O(1/m)$ , so that

$$\begin{aligned} \pi^R\left(\frac{v}{m}\right) &\propto \left[ \left(\frac{m}{v}\right)^2 \left(\frac{1}{m} + \frac{v(m-1)}{m^2} \sum_{i=1}^{n-1} \frac{1}{v+i}\right) + O(1) + \sum_{j=1}^{n-1} \frac{1}{(\frac{v}{m} + j)^2} O\left(\frac{1}{m}\right) - \sum_{i=0}^{n-1} \frac{m}{(v+i)^2} \right]^{1/2} \\ &= \left[ \sum_{i=1}^{n-1} \frac{1}{v+i} \left(\frac{(m-1)}{v} - \frac{m}{(v+i)}\right) + O(1) \right]^{1/2} \\ &= \left[ \frac{(m-1)}{v} \sum_{i=1}^{n-1} \frac{i}{(v+i)^2} + O(1) \right]^{1/2} \\ &= \sqrt{m-1} \left[ \frac{1}{v} \sum_{i=1}^{n-1} \frac{i}{(v+i)^2} + O\left(\frac{1}{m}\right) \right]^{1/2}. \end{aligned}$$

Combining this with (11), noting that  $v\Gamma(v) = \Gamma(v+1)$ , and letting  $m \rightarrow \infty$ , yields the result.  $\square$

It follows, of course, that  $a$  behaves like  $v/m$  for large  $m$ , where  $v$  has the distribution in (10). It is very interesting that this “large  $m$ ” behavior of the posterior depends on the data only through  $r_0$ , the number of nonzero cell observations.

If, in addition,  $n$  is moderately large (but much smaller than  $m$ ), we can explicitly study the behavior of the posterior mode of  $a$ .

**Proposition 3.5** *Suppose  $m \rightarrow \infty$ ,  $n \rightarrow \infty$ , and  $n/m \rightarrow 0$ . Then (10) has mode*

$$\hat{v} \approx \begin{cases} \frac{(r_0-1.5)}{\log(1+n/r_0)} & \text{if } \frac{r_0}{n} \rightarrow 0, \\ c^*n & \text{if } \frac{r_0}{n} \rightarrow c < 1, \end{cases}$$

where  $r_0$  is the number of nonzero  $x_i$ ,  $c^*$  is the solution to  $c^* \log(1 + \frac{1}{c^*}) = c$ , and  $f(n, m) \approx g(n, m)$  means  $f(n, m)/g(n, m) \rightarrow 1$ . The corresponding mode of the reference posterior for  $a$  is  $\hat{a}^R = \hat{v}/m$ .

*Proof.* Taking the log of (10) and differentiating with respect to  $v$  results in

$$\Psi'(v) = \frac{(r_0 - 1.5)}{v} - \sum_{i=1}^{n-1} \frac{1}{v+i} - \frac{\sum_{i=1}^{n-1} \frac{i}{(v+i)^3}}{\sum_{i=1}^{n-1} \frac{i}{(v+i)^2}}.$$

Note first that, as  $n$  grows, and if  $v$  also grows (no faster than  $n$ ), then

$$\sum_{i=1}^{n-1} \frac{1}{v+i} = \int_1^n \frac{1}{v+x} dx + O\left(\frac{1}{v+1}\right) + O\left(\frac{1}{n}\right) = \log\left(\frac{v+n}{v+1}\right) + O\left(\frac{1}{v+1}\right) + O\left(\frac{1}{n}\right).$$

Next,

$$\begin{aligned} \sum_{i=1}^{n-1} \frac{i}{(v+i)^3} &= \int_1^n \frac{x}{(v+x)^3} dx + O\left(\frac{1}{(v+1)^2}\right) + O\left(\frac{1}{n^2}\right) \\ &= \frac{1}{2} \left[ \frac{(v+2)}{(v+1)^2} - \frac{(v+2n)}{(v+n)^2} \right] + O\left(\frac{1}{(v+1)^2} + \frac{1}{n^2}\right) = O\left(\frac{1}{v+1}\right) + O\left(\frac{1}{n}\right), \end{aligned}$$

$$\begin{aligned} \sum_{i=1}^{n-1} \frac{i}{(v+i)^2} &= \int_1^n \frac{x}{(v+x)^2} dx + O\left(\frac{1}{(v+1)}\right) + O\left(\frac{1}{n}\right) \\ &= \frac{v(1+n)}{(v+1)(v+n)} + \log\left(\frac{v+n}{v+1}\right) + O\left(\frac{1}{v+1} + \frac{1}{n}\right) \geq \log 2, \end{aligned}$$

again using that  $v$  will not grow faster than  $n$ . Putting these together we have that

$$\Psi'(v) = \frac{(r_0 - 1.5)}{v} - \log\left(\frac{v+n}{v+1}\right) + O\left(\frac{1}{v+1}\right) + O\left(\frac{1}{n}\right).$$

*Case 1.*  $\frac{r_0}{n} \rightarrow c$ , for  $0 < c < 1$ . For this case, write  $v = c^*n/(1+\delta)$  for  $\delta$  small, and note that then

$$\Psi'(v) = \frac{c}{c^*}(1+o(1))(1+\delta) - \log\left(\frac{(c^*+1)}{c^*}\right) + o(1).$$

Since  $\frac{c}{c^*} - \log\left(\frac{(c^*+1)}{c^*}\right) = 0$ , it is clear that  $\delta$  can be appropriately chosen as  $o(1)$  to make the derivative zero.

*Case 2.*  $\frac{r_0}{n} \rightarrow 0$ . Now choose  $v = \frac{(r_0-1.5)}{(1+\delta)\log(1+n/r_0)}$  and note that  $\frac{v}{n} \rightarrow 0$ . It follows that

$$\log\left(1 + \frac{n}{r_0}\right) = (\log n - \log r_0 + o(1))(1+\delta) \quad \text{and} \quad \log\left(\frac{v+n}{v+1}\right) = [\log n - \log(v+1)](1+o(1)).$$

Consider first the case  $v \rightarrow \infty$ . Then

$$\log(v+1) = (1+o(1))(\log r_0 - \log \log(1+n/r_0)) = (1+o(1)) \log r_0,$$

so that

$$\Psi'(v) = (\log n - \log r_0 + o(1))(1 + \delta) - (\log n - \log r_0)(1 + o(1)) + o(1),$$

and it is clear that  $\delta$  can again be chosen  $o(1)$  to make this zero. Lastly, if  $v \leq K < \infty$ , then  $(\log r_0)/(\log n) = o(1)$ , so that  $\Psi'(v) = (\log n)(1 + o(1))(1 + \delta) - (\log n)(1 + o(1)) + o(1)$ , and  $\delta$  can again be chosen  $o(1)$  to make this zero, completing the proof.  $\square$

Table 1 gives the limiting behavior of  $\hat{v}$  for various behaviors of the number of nonzero cells,  $r_0$ . Only when  $r_0 = \log n$  does the posterior mode of  $a$  (i.e.,  $v/m$ ) equal  $1/m$ , the value selected by the reference distance method. Of course, this is not surprising; empirical Bayes is using a fit to the data to help select  $a$  whereas the reference distance method is pre-experimental.

$r_0$	$cn$ ( $0 < c < 1$ )	$n^b$ ( $0 < b < 1$ )	$(\log n)^b$	$\log n$	$O(1)$
$\hat{v}$	$c^*n$	$\frac{n^b}{(1-b)\log n}$	$(\log n)^{(b-1)}$	1	$O(1/\log n)$

Table 1: *The limiting behavior of  $\hat{v}$  as  $n \rightarrow \infty$ , for various limiting behaviors of  $r_0$ , the number of non-zero cells.*

### 3.2 Multi-normal means

Let  $x_i$  be independent normal with mean  $\mu_i$  and variance 1, for  $i = 1 \dots, m$ . We are interested in all the  $\mu_i$  and in  $|\boldsymbol{\mu}|^2 = \mu_1^2 + \dots + \mu_m^2$ .

The natural hierarchical prior modeling approach is to assume that  $\mu_i \stackrel{iid}{\sim} N(\mu_i | 0, \tau)$ . Then, marginally, the  $x_i$  are iid  $N(x_i | 0, \sqrt{1 + \tau^2})$  and the reference (Jeffreys) prior for  $\tau^2$  in this marginal model is

$$\pi^R(\tau^2) \propto (1 + \tau^2)^{-1}.$$

The hierarchical prior for  $\boldsymbol{\mu}$  (and recommended overall prior) is then

$$\pi^o(\boldsymbol{\mu}) = \int_0^\infty \frac{1}{(2\pi\tau^2)^{m/2}} \exp\left(-\frac{|\boldsymbol{\mu}|^2}{2\tau^2}\right) \frac{1}{1 + \tau^2} d\tau^2. \quad (12)$$

This prior is arguably reasonable from a marginal reference prior perspective. For the individual  $\mu_i$ , it is a shrinkage prior known to result in Stein-like shrinkage estimates of the form

$$\hat{\mu}_i = \left(1 - \frac{r(|\boldsymbol{x}|)}{|\boldsymbol{x}|^2}\right) x_i,$$

with  $r(\cdot) \approx p$  for large arguments. Such shrinkage estimates are often viewed as actually being superior to the reference posterior mean, which is just  $x_i$  itself. The reference prior



when  $|\mu|$  is the parameter of interest is

$$\pi_{|\mu|}(\boldsymbol{\mu}) \propto \frac{1}{|\boldsymbol{\mu}|^{m-1}} \propto \int_0^\infty \frac{1}{(2\pi\tau^2)^{m/2}} \exp\left(-\frac{|\boldsymbol{\mu}|^2}{2\tau^2}\right) \frac{1}{\tau} d\tau^2, \quad (13)$$

which is clearly very similar to (12). Thus the hierarchical prior appears to be quite satisfactory in terms of its marginal posterior behavior for any of the parameters of interest. Of course, the same could be said for the single reference prior in (13); thus here is a case where one of the reference priors would be fine for all parameters of interest, and averaging among reference priors would not work.

Computation with the reference prior in (13) can be done by a simple Gibbs sampler. Computation with the hierarchical prior in (12) is almost as simple, with the Gibbs step for  $\tau^2$  being replaced by the rejection step:

*Step 1.* Propose  $\tau^2$  from the inverse gamma density proportional to

$$\frac{1}{(\tau^2)^{(1+m/2)}} \exp\left(-\frac{|\boldsymbol{\mu}|^2}{2\tau^2}\right),$$

*Step 2.* Accept the result with probability  $\tau^2/(\tau^2 + \sigma^2/n)$  (or else propose again).

### 3.3 Bivariate normal problem

Earlier for the bivariate normal problem, we only considered the two right-Haar priors. More generally, there is a continuum of right-Haar priors given as follows. Define an orthogonal matrix by

$$\mathbf{\Gamma} = \begin{pmatrix} \cos(\beta) & -\sin(\beta) \\ \sin(\beta) & \cos(\beta) \end{pmatrix}$$

where  $-\pi/2 < \beta \leq \pi/2$ . Then it is straightforward to see that the right-Haar prior based on the transformed data  $\mathbf{\Gamma}\mathbf{X}$  is

$$\pi(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho \mid \beta) = \frac{\sin^2(\beta) \sigma_1^2 + \cos^2(\beta) \sigma_2^2 + 2 \sin(\beta) \cos(\beta) \rho \sigma_1 \sigma_2}{\sigma_1^2 \sigma_2^2 (1 - \rho^2)}.$$

We thus have a class of priors indexed by a hyperparameter  $\beta$ , and can apply the hierarchical approach to obtain an overall prior. The natural prior distribution on  $\beta$  is the (proper) uniform distribution (being uniform over the set of rotations is natural.) The resulting joint prior is

$$\pi^o(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho) = \frac{1}{\pi} \int_{-\pi/2}^{\pi/2} \pi(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho \mid \beta) d\beta,$$

which equals the prior  $\pi^A$  in (1), since

$$\int_{-\pi/2}^{\pi/2} \sin(\beta) \cos(\beta) d\beta = 0, \quad \int_{-\pi/2}^{\pi/2} \sin^2(\beta) d\beta = \int_{-\pi/2}^{\pi/2} \cos^2(\beta) d\beta = \text{constant}.$$

Thus the overall prior obtained by the hierarchical approach is the same prior as obtained by just averaging the two reference priors. It was stated there that this prior is inferior as an overall prior to either reference prior individually, so the hierarchical approach has failed.

**Empirical hierarchical approach:** Instead of integrating out over  $\beta$ , one could find the empirical Bayes estimate  $\hat{\beta}$  and use  $\pi(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho \mid \hat{\beta})$  as the overall prior. This was shown in Sun and Berger (2007) to result in a terrible overall prior, much worse than either the individual reference priors, or even  $\pi^A$  in (1).

## 4 Discussion

The reference distance approach to developing an overall prior is natural, and seems to work well when the reference priors themselves are proper. It also appears to be possible to implement the approach in the case where the reference priors are improper, by operating on suitable large compact sets and showing that the result is not sensitive to the choice of compact set.

The hierarchical approach seems excellent (as usual), and can certainly be recommended if one can find a natural hierarchical structure based on a class of proper priors. In particular, the overall prior obtained for the multi-normal mean problem seems fine, and the recommended hierarchical prior for the contingency table situation is very interesting, and seems to have interesting adaptations to sparsity; the same can be said for its empirical Bayes implementation. In contrast, the hierarchical and empirical Bayes implementations were very unsatisfactory for the bivariate normal problem, when based on the class of right-Haar priors, even though the hyperprior was proper. This is a clear warning to use the hierarchical or empirical Bayes approach only with a base class of proper priors.

The failure of arithmetic prior averaging in the bivariate normal problem was also dramatic; the initial averaging of two right-Haar priors gave an inferior result, which was duplicated by the continuous average over all right-Haar priors. Curiously in this example, the geometric average of the two right-Haar improper priors seems to be reasonable, suggesting that, if averaging of improper priors is to be done, the geometric average should be used.

## Acknowledgement

Berger's work was supported by NSF Grants DMS-0757549-001 and DMS-1007773. Sun's work was supported by NSF grants DMS-1007874 and SES-1260806. The research is also supported by Chinese 111 Project B14019.

## References

- 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). 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. and Bernardo, J. M. (1992b). Ordered group reference priors, with applications to multinomial problems. *Biometrika* **79**, 25–37.
- Berger, J., 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. (2012). Objective priors for discrete parameter spaces. *J. Amer. Statist. Assoc.* **107**, 636–648.
- Berger and Sun, D. (2008). Objective priors for the bivariate normal model *Ann. Statist.* **36**, 963–982.
- Bernardo, J. M. (1979). Reference posterior distributions for Bayesian inference. *J. Roy. Statist. Soc. B* **41**, 113–147 (with discussion).
- Bernardo, J. M. (2005). Reference analysis. *Bayesian Thinking: Modeling and Computation, Handbook of Statistics* **25** (D. K. Dey and C. R. Rao, eds.) Amsterdam: Elsevier, 17–90.
- 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. (2006). Intrinsic point estimation of the normal variance. *Bayesian Statistics and its Applications*. (S. K. Upadhyay, U. Singh and D. K. Dey, eds.) New Delhi: Anamaya Pub, 110–121.
- Bernardo, J. M. and Smith, A. F. M. (1994). *Bayesian Theory*. Chichester: Wiley.
- Clarke, B. and Barron, A. (1994). Jeffreys’ prior is the reference prior under entropy loss. *J. Statist. Planning and Inference* **41**, 37–60.
- Clarke, B. and Yuan A. (2004). Partial information reference priors: derivation and interpretations. *J. Statist. Planning and Inference* **123**, 313–345.
- Consonni, G., Veronese, P. and Gutiérrez-Peña E. (2004). Reference priors for exponential families with simple quadratic variance function. *J. Multivariate Analysis* **88**, 335–364.
- Datta, G. S. and Ghosh, J. K. (1995a). On priors providing frequentist validity for Bayesian inference. *Biometrika* **82**, 37–45.

- Datta, G. S. and Ghosh, J. K. (1995b). Noninformative priors for maximal invariant parameter in group models. *Test* **4**, 95–114.
- Datta, G. S. and Ghosh, M. (1996). On the invariance of noninformative priors. *Ann. Statist.* **24**, 141–159.
- Datta, G. S., Mukerjee, R., Ghosh, M. and Sweeting, T. J. (2000). Bayesian prediction with approximate frequentist validity. *Ann. Statist.* **28** 1414–1426.
- De Santis, F., Morteo, J. and Nardi, A. (2001). Jeffreys priors for survival models with censored data. *J. Statist. Planning and Inference* **99**, 193–209.
- De Santis, F. (2006). Power priors and their use in clinical trials. *The American Statistician*, **60**. 122–129.
- Ghosh, J. K. and Ramamoorthi, R. V. (2003). *Bayesian Nonparametrics*. New York: Springer
- Ghosh, M., Mergel, V., and Liu, R. (2011). A general divergence criterion for prior selection. *Ann. Inst. Statist. Math.* **60**, 43–58.
- Ghosh, M. (2011). Objective priors: An introduction for frequentists. *Statistical Science* **26**, 187–202.
- Gilks, W.R. and Wild, P. (1992). Adaptive rejection sampling for Gibbs sampling. *Appl. Statist.* **41**, 337–348.
- Hartigan, J. A. (1964). Invariant prior distributions. *Ann. Math. Statist.* **35**, 836–845.
- Jeffreys, H. (1946). An invariant form for the prior probability in estimation problems. *Proc. Roy. Soc. A* **186**, 453–461.
- Jeffreys, H. (1961). *Theory of Probability* (3rd edition). Oxford: Oxford University Press.
- Kass, R. E. and Wasserman, L. (1996). The selection of prior distributions by formal rules. *J. Amer. Statist. Assoc.* **91**, 1343–1370.
- Kullback, S. and R. A. Leibler, R. A. (1951). On information and sufficiency. *Ann. Math. Statist.* **22**, 79–86
- Laplace, P. S. (1812). *Thorie Analitique des Probabilits*. Paris: Courcier. Reprinted as *Oeuvres Compltes de Laplace* **7**, 1878–1912. Paris: Gauthier-Villars.
- Liseo, B. (1993). Elimination of nuisance parameters with reference priors. *Biometrika*, **80**, 295–304.
- Liseo, B, and Loperfido, N, (2006). A note on reference priors for the scalar skew-normal distribution, *J. Statist. Planning and Inference* **136**, 373–389.
- Sivaganesan, S. (1994). Discussion to "An Overview of Bayesian Robustness" by J. Berger. *Test*, **3**, 116–120.

- Sivaganesan, S., Laud, P., Mueller, P. (2011). A Bayesian subgroup analysis using zero-inflated Polya-urn scheme. *Sociol. Methodology* **30**, 312–323.
- Stein, C. (1959). An example of wide discrepancy between fiducial and confidence intervals. *Ann. Math. Statist.* **30**, 877–880.
- Stein, C. (1964). Inadmissibility of the usual estimator for the variance of a normal distribution with unknown mean. *Ann. Inst. Statist. Math.* **16**, 155–160.
- Stone, M. and Dawid, A. P. (1972). Un-Bayesian implications of improper Bayesian inference in routine statistical problems. *Biometrika* **59**, 269–375.
- Sun, D. and Berger, J. O. (2007). Objective Bayesian analysis for the multivariate normal model. *Bayesian Statistics 8* (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, 525–562 (with discussion).
- Walker, S. G. and Gutiérrez-Peña, E. (2011). A decision-theoretical view of default priors *Theory and Decision* **70**, 1–11
- Yang, R. and Berger, J. O. (1994). Estimation of a covariance matrix using the reference prior. *Ann. Statist.* **22**, 1195–2111.