



## Robust Bayes and Empirical Bayes Analysis with # -Contaminated Priors

James Berger; L. Mark Berliner

*The Annals of Statistics*, Vol. 14, No. 2. (Jun., 1986), pp. 461-486.

Stable URL:

<http://links.jstor.org/sici?sici=0090-5364%28198606%2914%3A2%3C461%3ARBAEBA%3E2.0.CO%3B2-J>

*The Annals of Statistics* is currently published by Institute of Mathematical Statistics.

---

Your use of the JSTOR archive indicates your acceptance of JSTOR's Terms and Conditions of Use, available at <http://www.jstor.org/about/terms.html>. JSTOR's Terms and Conditions of Use provides, in part, that unless you have obtained prior permission, you may not download an entire issue of a journal or multiple copies of articles, and you may use content in the JSTOR archive only for your personal, non-commercial use.

Please contact the publisher regarding any further use of this work. Publisher contact information may be obtained at <http://www.jstor.org/journals/ims.html>.

Each copy of any part of a JSTOR transmission must contain the same copyright notice that appears on the screen or printed page of such transmission.

---

JSTOR is an independent not-for-profit organization dedicated to and preserving a digital archive of scholarly journals. For more information regarding JSTOR, please contact [support@jstor.org](mailto:support@jstor.org).

## ROBUST BAYES AND EMPIRICAL BAYES ANALYSIS WITH $\varepsilon$ -CONTAMINATED PRIORS

BY JAMES BERGER<sup>1</sup> AND L. MARK BERLINER<sup>2</sup>

*Purdue University and Ohio State University*

For Bayesian analysis, an attractive method of modelling uncertainty in the prior distribution is through use of  $\varepsilon$ -contamination classes, i.e., classes of distributions which have the form  $\pi = (1 - \varepsilon)\pi_0 + \varepsilon q$ ,  $\pi_0$  being the base elicited prior,  $q$  being a "contamination," and  $\varepsilon$  reflecting the amount of error in  $\pi_0$  that is deemed possible. Classes of contaminations that are considered include (i) all possible contaminations, (ii) all symmetric, unimodal contaminations, and (iii) all contaminations such that  $\pi$  is unimodal.

Two issues in robust Bayesian analysis are studied. The first is that of determining the range of posterior probabilities of a set as  $\pi$  ranges over the  $\varepsilon$ -contamination class. The second, more extensively studied, issue is that of selecting, in a data dependent fashion, a "good" prior distribution (the Type-II maximum likelihood prior) from the  $\varepsilon$ -contamination class, and using this prior in the subsequent analysis. Relationships and applications to empirical Bayes analysis are also discussed.

**1. Introduction.** The most frequent criticism of subjective Bayesian analysis is that it supposedly presumes an ability to completely and accurately quantify subjective information in terms of a single prior distribution. However, there has long existed [at least since Good (1950)] a *robust Bayesian viewpoint* which assumes only that subjective information can be quantified in terms of a class  $\Gamma$  of possible distributions. The goal is then to make inferences or decisions which are *robust* over  $\Gamma$ , i.e., which are relatively insensitive (or at least are satisfactory) to deviations as the prior distribution varies over  $\Gamma$ . We will not consider the philosophical or pragmatic reasons for adopting this viewpoint. Such a discussion, along with a review of the area, may be found in Berger (1984). (We also do not mean to imply that the single prior Bayesian approach is necessarily bad; it usually works very well.) Related to this are various forms of empirical Bayes analysis [cf. Morris (1983) for discussion and review], in which the prior distribution is also assumed to belong to some class  $\Gamma$  of distributions. Indeed, Section 5 considers some familiar empirical Bayes problem from our perspective. Also, see Berger and Berliner (1984).

Before discussing implementation of the robust Bayesian viewpoint, some notation is helpful. Let  $X$  denote the observable random variable (or vector), which will (for simplicity) be assumed to have a density  $f(x|\theta)$  (w.r.t. some

---

Received November 1983; revised October 1985.

<sup>1</sup>Research supported by the National Science Foundation under Grants MCS-8101670A1 and DMS-8401996.

<sup>2</sup>Research supported by the Office of Naval Research under Contract N00014-84-K-0422.

AMS 1980 *subject classifications*. Primary 62A15; secondary 62F15.

*Key words and phrases*. Robust Bayes, empirical Bayes, classes of priors,  $\varepsilon$ -contamination, type II maximum likelihood, hierarchical priors.

measure), where  $\theta$  is an unknown parameter lying in a parameter space  $\Theta$ . A prior distribution on  $\Theta$  will be denoted by  $\pi$  (later, in examples,  $\pi$  will be used to denote either a prior or its corresponding density), and the resulting marginal density of  $X$  is given by

$$m(x|\pi) = E^\pi f(x|\theta) = \int_{\Theta} f(x|\theta)\pi(d\theta).$$

The posterior distribution of  $\theta$  given  $x$  (assuming it exists) will be denoted by  $\pi(\cdot|x)$  and, in nice situations, is defined by

$$\pi(d\theta|x) = f(x|\theta)\pi(d\theta)/m(x|\pi).$$

Finally, let  $\mathcal{P}$  denote the space of all probability distributions on  $\Theta$ .

The class,  $\Gamma$ , of prior distributions to be considered in this paper, is the  $\varepsilon$ -contamination class; namely,

$$(1.1) \quad \Gamma = \{\pi: \pi = (1 - \varepsilon)\pi_0 + \varepsilon q, q \in \mathcal{Q}\},$$

where  $0 \leq \varepsilon \leq 1$  is given,  $\pi_0$  is a particular prior distribution, and  $\mathcal{Q}$  is some subset of  $\mathcal{P}$ . There are several reasons for consideration of this class. First, and foremost, it is a sensible class to consider in light of the prior elicitation process. The extensive and rapidly developing methodology on prior elicitation [cf. Kadane et al. (1980)] makes specification of an initial believable prior,  $\pi_0$ , an attractive starting point. However, in determining  $\pi_0$  sensibly, one will make *probability* judgements about subsets of  $\Theta$ , judgements which could be in error by some amount  $\varepsilon$ . Stated another way, further reflection might lead to alterations of probability judgements by an amount  $\varepsilon$ . Hence, possible priors involving such alterations should be included in  $\Gamma$ .

Many classes of priors which have been considered are not sensible from the above viewpoint. For instance, classes of priors involving restrictions on moments force severe restrictions on the allowable prior tails. This makes little sense from the elicitation viewpoint, since the tails of a prior involve very small probabilities and are, therefore, nearly impossible to determine. Similarly, classes of conjugate priors are too limited, particularly in their inflexible tail behavior [cf. Berger (1985)].

Two other major reasons for choosing  $\Gamma$  as in (1.1) are (i) such  $\Gamma$  are (as we shall see) surprisingly easy to work with; and (ii) such  $\Gamma$  are very flexible through choice of  $\mathcal{Q}$ . In this paper we will restrict consideration to four interesting choices of  $\mathcal{Q}$ . First, in Section 2 the choice  $\mathcal{Q} = \mathcal{P}$  (all distributions) will be considered. This choice is easy to work with and is, in some sense, conservative. In Section 3 we consider the class,  $\mathcal{Q}$ , of all contaminations which are symmetric and unimodal. This is again very easy to work with. In Section 4 we consider the class of all contaminations such that *the resulting*  $\pi$  is unimodal (assuming that  $\pi_0$  is unimodal). It came as a great surprise to us that such a complicated class could be worked with and provide reasonably simple answers. Finally, in Section 5 we consider  $\mathcal{Q}$  that are mixtures of various classes. The purpose of the section is to show how easily mixed contaminations can be dealt with and also to apply the methodology in some typical empirical Bayes situations.

Other articles that have used  $\varepsilon$ -contamination classes of priors include Schneeweiss (1964), Blum and Rosenblatt (1967), Huber (1973), Marazzi (1985), Bickel (1984), and Berger (1982, 1984). Except for Huber (1973), these articles work within the frequentist Bayesian framework, whereas our approach will be almost entirely conditional Bayesian. Huber (1973) is discussed below and in Section 2.4. There is a substantial literature working with other types of classes of priors [cf. Leamer (1978) and DeRobertis and Hartigan (1981)], and with the very related idea of “upper” and “lower” probabilities. Berger (1984) contains considerable review and discussion of this literature. We strongly prefer the class in (1.1) for intuitive content and ease of analysis.

The ideal analysis, to a robust Bayesian, is one in which it can be shown that the inference or decision to be made is essentially the same for any prior in  $\Gamma$ . [Indeed, it can be argued—see Berger (1984, 1985)—that this is the only way in which a statistical conclusion can claim to be ultimately sound.] What is needed, to provide such conclusions, is essentially the ability to find minimums and maximums of criterion functions as  $\pi$  ranges over  $\Gamma$ . We illustrate this approach in Section 2.4, where, for  $\mathcal{Q} = \{\text{all distributions}\}$ , the range of posterior probabilities of a (fixed) set  $C$  is given [essentially following Huber (1973)]. This allows finding the range of posterior probabilities of confidence sets and the range of posterior probabilities of hypotheses, for such  $\Gamma$ .

Unfortunately, there are certain inadequacies in assuming that  $\mathcal{Q} = \{\text{all distributions}\}$  (see Section 2.3), and attempting the above program with more reasonable  $\Gamma$  (such as that in Section 4) becomes more difficult. A number of alternative approaches to the problem of dealing with classes of priors have thus been proposed, essentially leading to the choice of a single “robust” prior, decision, or inference. Berger (1984, 1985) discusses various of these methods, including the appealing technique of putting a prior distribution on  $\Gamma$ . [Such a prior is called a hyperprior or a Type II probability distribution by Good (1965, 1980) and a second-stage prior in certain situations by Lindley and Smith (1972).] Of course, this corresponds to using a certain single prior (the “average” over  $\Gamma$ ), but one would suspect that the resulting Bayes rule would be quite robust with respect to  $\Gamma$ . The difficulty in doing this is mainly technical: it is essentially impossible to put a reasonable prior on complicated  $\Gamma$ , such as those in Sections 2–4, and carry out the Bayesian calculations. Note also that, ideally, most of the prior information available will have been exhausted in constructing  $\Gamma$ . Hence, any prior distribution placed on  $\Gamma$  will to a large extent, be arbitrary.

Instead, we will consider the simplest and most commonly used method of selecting a hopefully robust prior in  $\Gamma$ , namely choice of that prior  $\pi$  which maximizes the marginal  $m(x|\pi)$  over  $\Gamma$ . This process is called Type II maximum likelihood by Good (1965). For  $\pi = (1 - \varepsilon)\pi_0 + \varepsilon q$ ,  $q \in \mathcal{Q}$ , maximizing  $m(x|\pi) = (1 - \varepsilon)m(x|\pi_0) + \varepsilon m(x|q)$  over  $\pi$  is clearly done by maximizing  $m(x|q)$  over  $q$ . Assuming that the maximum of  $m(x|q)$  is attained at (a unique)  $\hat{q} \in \mathcal{Q}$ , we will then suggest formally using the estimated prior  $\hat{\pi}$ , given by

$$(1.2) \quad \hat{\pi} = (1 - \varepsilon)\pi_0 + \varepsilon\hat{q}.$$

(Of course,  $\hat{\pi}$  thus depends on  $x$ .) Throughout the paper,  $\hat{\pi}$  will be called the

ML-II prior. Also, any quantities derived from  $\hat{\pi}$  will appear with the modifier "ML-II" for clarity.

Choosing a prior with the help of the data always engenders controversy. Several justifications for doing so can be given, however. First, if  $m(x|\pi)$  is "small," it is simply unlikely that such a  $\pi$  could be "true," and hence worrying about such  $\pi$  is counterproductive. Recall that (supposedly) all  $\pi \in \Gamma$  are deemed to be reasonable representations of prior beliefs, so  $\hat{\pi}$  is simply the prior which is most plausible, in light of prior opinions and the data. A more formal way of saying this is that, if all  $\pi \in \Gamma$  are roughly equally likely a priori, then  $\hat{\pi}$  is the "posterior mode" of the "uniform" distribution on  $\Gamma$ , and might often be expected to yield a posterior distribution that is close to the true posterior distribution for such a "uniform" distribution on  $\Gamma$ .

The preceding argument for  $\hat{\pi}$  is, of course, nonrigorous, and the ultimate justification for proceeding in this way is simply that it can give reasonable answers. Of course, there is already substantial evidence in the literature attesting to the success of the method, both in the Bayesian literature [cf. Jeffreys (1961), Good (1965, 1980), Box and Tiao (1973), Bishop, Fienberg, and Holland (1975), and Zellner (1985)], and in the empirical Bayesian literature [cf. Maritz (1970) and Morris (1983)]. Indeed, note that the "standard" empirical Bayes methodology is to choose  $\Gamma$  to be a class of conjugate priors and then to estimate the "hyperparameters" of the prior by maximizing  $m(x|\pi)$ , yielding  $\hat{\pi}$ . Also the related use of the marginal in Bayesian model robustness investigations is well established [cf. Box and Tiao (1973), Dempster (1975), and Box (1980)]. When all is said and done, however, we recognize that the ML-II technique is not foolproof and can produce bad answers, particularly when  $\Gamma$  includes unreasonable distributions. (The basic problem with the ML-II technique is that ensuing calculations of variability do not take into account the "error" of the ML-II estimation; see Section 2.3 for an extreme example of this problem.) In Section 6 we give a general discussion of the success of the method for the situations discussed in the paper.

We conclude this section with useful formulas and notation. For priors of the form

$$(1.3) \quad \pi(d\theta) = (1 - \varepsilon)\pi_0(d\theta) + \varepsilon q(d\theta),$$

computations give [assuming the existence of the posterior distributions  $\pi_0(d\theta|x)$  and  $q(d\theta|x)$ ]

$$(1.4) \quad m(x|\pi) = (1 - \varepsilon)m(x|\pi_0) + \varepsilon m(x|q)$$

and

$$(1.5) \quad \pi(d\theta|x) = \lambda(x)\pi_0(d\theta|x) + (1 - \lambda(x))q(d\theta|x),$$

where  $\lambda(x) \in [0, 1]$  is given by

$$(1.6) \quad \lambda(x) = (1 - \varepsilon)m(x|\pi_0)/m(x|\pi).$$

Furthermore, the posterior mean,  $\delta^\pi$ , and posterior variance,  $V^\pi$ , can be written

(assuming they exist) as

$$(1.7) \quad \delta^\pi(x) = \lambda(x)\delta^{\pi_0}(x) + (1 - \lambda(x))\delta^q(x)$$

and

$$(1.8) \quad V^\pi(x) = \lambda(x)V^{\pi_0}(x) + (1 - \lambda(x))V^q(x) + \lambda(x)(1 - \lambda(x))(\delta^{\pi_0}(x) - \delta^q(x))^2.$$

Part of the appeal of the  $\epsilon$ -contamination class,  $\Gamma$ , is the simplicity of these formulas.

**2. Analysis for arbitrary contaminations.** A natural suggestion for a class of contaminations of a fixed, elicited prior  $\pi_0$  is the class of all possible contaminations. In this section we will examine inferences, including point estimation, testing, and credible regions, for such a class, i.e., for

$$(2.1) \quad \Gamma = \{ \pi: \pi = (1 - \epsilon)\pi_0 + \epsilon q, q \in \mathcal{P} \}.$$

In a number of respects this is too *large* a class of priors, including many priors that are unreasonable. And it will be seen that this can lead to serious difficulties in some situations (although for certain purposes no problems are encountered). We give a fairly detailed analysis of this situation because its relative simplicity allows easy comprehension of important concepts (*including* the difficulties of using too large a  $\Gamma$ ), and because *some* useful robustness results do emerge. All proofs are easy and are omitted.

*2.1. The ML-II prior and posterior.* For  $\Gamma$  defined as in (2.1), the ML-II prior and corresponding posterior are as follows.

**THEOREM 2.1.** *Assume  $X$  has a density  $f(x|\theta)$  w.r.t. some dominating measure on the sample space of  $X$ . Assume that the usual maximum likelihood estimator for  $\theta$ , say  $\hat{\theta}(x)$ , exists and is unique. For  $\Gamma$  defined as in (2.1), the ML-II prior is given by*

$$(2.2) \quad \hat{\pi}(\cdot) = (1 - \epsilon)\pi_0(\cdot) + \epsilon\hat{q}_x(\cdot),$$

where  $\hat{q}_x$  assigns probability one to the point  $\theta = \hat{\theta}(x)$ . The ML-II posterior is given by

$$(2.3) \quad \hat{\pi}(\cdot|x) = \hat{\lambda}(x)\pi_0(\cdot|x) + (1 - \hat{\lambda}(x))\hat{q}_x(\cdot),$$

where

$$(2.4) \quad \hat{\lambda}(x) = (1 - \epsilon)m(x|\pi_0) / [(1 - \epsilon)m(x|\pi_0) + \epsilon f(x|\hat{\theta}(x))].$$

*2.2. The ML-II posterior mean.* Under the assumptions of Theorem 2.1, the ML-II posterior mean of  $\theta$  is given by [see (1.7)]

$$(2.5) \quad \delta^{\hat{\pi}}(x) = \hat{\lambda}(x)\delta^{\pi_0}(x) + (1 - \hat{\lambda}(x))\hat{\theta}(x).$$

As an estimator of  $\theta$ ,  $\delta^{\hat{\pi}}$  is intuitively appealing, in that it is a reasonable data dependent mixture of  $\delta^{\pi_0}$  and  $\hat{\theta}$ . When the data are consistent with  $\pi_0$ ,  $m(x|\pi_0)$

will be reasonably large and  $\hat{\lambda}(x)$  close to one (for small  $\varepsilon$ ), so that  $\delta^{\hat{\pi}}$  will essentially equal  $\delta^{\pi_0}$ . When the data and  $\pi_0$  are not compatible, however,  $m(x|\pi_0)$  will be small and  $\hat{\lambda}(x)$  near zero;  $\delta^{\hat{\pi}}$  will then be approximately equal to the m.l.e.  $\hat{\theta}$ .

The following example presents  $\delta^{\hat{\pi}}$  in an important situation. Some properties of the estimator are discussed which give a degree of “outside validation” to the estimator.

**EXAMPLE 1.** Let  $X = (X_1, \dots, X_p)^t \sim \mathcal{N}_p(\theta, \sigma^2 I_p)$ , where  $\theta = (\theta_1, \dots, \theta_p)^t$  is unknown and  $\sigma^2$  is known. Suppose the elicited prior,  $\pi_0$ , for  $\theta$  is  $\mathcal{N}_p(\mu, \tau^2 I_p)$ . (Thus  $\mu$  and  $\tau^2$  are specified.) Since the usual maximum likelihood estimator of  $\theta$  is  $\hat{\theta}(x) = x$ , and  $\delta^{\pi_0}(x) = x - [\sigma^2/(\sigma^2 + \tau^2)](x - \mu)$ , formula (2.5) reduces to

$$\delta^{\hat{\pi}}(x) = (1 - \hat{\lambda}(x)\sigma^2/(\sigma^2 + \tau^2))(x - \mu) + \mu,$$

where

$$\hat{\lambda}(x) = \left[1 + (\varepsilon/(1 - \varepsilon))(1 + \tau^2/\sigma^2)^{p/2} \exp\{|x - \mu|^2/2(\sigma^2 + \tau^2)\}\right]^{-1}.$$

Note that  $\hat{\lambda}$  goes to 0 exponentially fast in  $|x - \mu|^2$ , so that  $\delta^{\hat{\pi}}(x) \rightarrow x$  quite rapidly as  $|x - \mu|^2$  gets large. Because of this, one might conjecture that the estimator is minimax, in a frequentist decision-theoretic sense under, say, quadratic loss. Unfortunately, this turns out not to be the case, although the deviation from minimaxity is usually fairly slight. It is also interesting to note that  $\delta^{\hat{\pi}}$  happens to coincide with the generalized Bayes estimator corresponding to the formal prior

$$\rho(d\theta) = (1 - \varepsilon)\pi_0(d\theta) + \varepsilon\rho_0(d\theta),$$

where  $\rho_0(d\theta) = (2\pi\sigma^2)^{p/2} d\theta$ . Note that priors of a similar form were considered by, for instance, Leonard (1974). The development here can be viewed as proposing a reasonable method for choosing the relative weights of  $\pi_0$  and  $(d\theta)$ .

*2.3. The ML-II posterior variance.* To determine the estimation error in using  $\delta^{\hat{\pi}}$ , it is natural to look at the posterior variance,  $V^{\hat{\pi}}$ . From (1.8) it follows that

$$V^{\hat{\pi}}(x) = \hat{\lambda}(x)\left[V^{\pi_0}(x) + (1 - \hat{\lambda}(x))(\delta^{\pi_0}(x) - \hat{\theta}(x))^2\right].$$

It will typically be the case (as in Example 1) that, as  $\hat{\lambda}(x) \rightarrow 0$ ,  $V^{\hat{\pi}}(x)$  will also go to zero. Indeed,  $\hat{\pi}$  will usually “converge” to a point mass at  $\hat{\theta}(x)$ . This is clearly inappropriate; although data incompatible with  $\pi_0$  can be cause for preference of  $\hat{\theta}(x)$  to  $\delta^{\pi_0}(x)$ , it does not cause one to think that  $\theta$  equals  $\hat{\theta}(x)$  exactly.

The trouble here is caused by the fact that  $\Gamma$  contains unrealistic distributions. We may feel that  $\pi_0$  could be in error, but surely a point mass at  $\hat{\theta}(x)$  (when far from the center of  $\pi_0$ ) is not usually a reasonable contamination to expect a priori. Working with  $\Gamma$  as in Sections 3 and 4, which do not allow such implausible contaminations, substantially alleviates this problem. (See also Section 6.)

2.4. *Robustness as  $\pi$  ranges over  $\Gamma$ .* As mentioned in the introduction, the ideal goal for a robustness study would be to show that a decision or inference being contemplated is satisfactory for all  $\pi \in \Gamma$ . When  $\mathcal{Q}$  is the class of all distributions, it often becomes feasible to check this. The basic tool is the following result of Huber (1973), concerning the range of posterior probabilities of a set.

**THEOREM 2.2** [Huber (1973)]. *Suppose  $\mathcal{Q} = \mathcal{P}$ . Let  $C$  be a measurable subset of  $\Theta$ , and define  $\beta_0$  to be the posterior probability of  $C$  under  $\pi_0$ , i.e.,*

$$\beta_0 = P^{\pi_0}(\theta \in C|X = x).$$

Then

$$(2.6) \quad \inf_{\pi \in \Gamma} P^\pi(\theta \in C|X = x) = \beta_0 \left\{ 1 + \frac{\varepsilon \sup_{\theta \notin C} f(x|\theta)}{(1 - \varepsilon)m(x|\pi_0)} \right\}^{-1},$$

and

$$(2.7) \quad \sup_{\pi \in \Gamma} P^\pi(\theta \in C|X = x) = \frac{(1 - \varepsilon)m(x|\pi_0)\beta_0 + \varepsilon \sup_{\theta \in C} f(x|\theta)}{(1 - \varepsilon)m(x|\pi_0) + \varepsilon \sup_{\theta \in C} f(x|\theta)}.$$

**EXAMPLE 2.** Assume that  $X \sim \mathcal{N}(\theta, \sigma^2)$ ,  $\sigma^2$  known, and that  $\pi_0$  is  $\mathcal{N}(\mu, \tau^2)$ . It is well known that  $\pi_0(d\theta|x)$  is  $\mathcal{N}(\delta(x), V^2)$ , where

$$\delta(x) = x - (\sigma^2/(\sigma^2 + \tau^2))(x - \mu), \quad V^2 = \sigma^2\tau^2/(\sigma^2 + \tau^2).$$

The usual  $100(1 - \alpha)\%$  Bayes credible region for  $\theta$  is

$$C = \{\theta: \delta(x) - K < \theta < \delta(x) + K\},$$

where  $K = z_{\alpha/2}V$ ,  $z_{\alpha/2}$  being the  $100(1 - \alpha/2)$  upper percentile of the standard normal distribution.

To investigate the robustness of  $C$ , we use (2.6) of Theorem 2.2. Note that

$$\sup_{\theta \notin C} f(x|\theta) = \begin{cases} (2\pi\sigma^2)^{-1/2} & \text{if } x \notin C, \\ (2\pi\sigma^2)^{-1/2} \exp\left\{-\frac{1}{2\sigma^2}(|x - \delta(x)| - K)^2\right\} & \text{if } x \in C. \end{cases}$$

Thus (2.6) becomes, for  $x \notin C$ ,

$$\inf_{\pi \in \Gamma} P^\pi(\theta \in C|X = x) = (1 - \alpha) \left\{ 1 + \frac{\varepsilon(1 + \tau^2/\sigma^2)^{1/2}}{(1 - \varepsilon)} \exp\left[\frac{(x - \mu)^2}{2(\sigma^2 + \tau^2)}\right] \right\}^{-1},$$

and, for  $x \in C$ ,

$$\begin{aligned} \inf_{\pi \in \Gamma} P(\theta \in C|X = x) &= (1 - \alpha) \left\{ 1 + \frac{\varepsilon(1 + \tau^2/\sigma^2)^{1/2}}{(1 - \varepsilon)} \right. \\ &\quad \left. \times \exp\left[\frac{(x - \mu)^2 - (|x - \mu|V/\tau - z_{\alpha/2}\tau)^2}{2(\sigma^2 + \tau^2)}\right] \right\}^{-1}. \end{aligned}$$



As a concrete example, suppose that  $\sigma^2 = 1$ ,  $\tau^2 = 2$ ,  $\mu = 0$ , and  $\varepsilon = 0.2$ . First, suppose  $x = 0.5$  is observed. Then the usual 95% Bayes credible interval for  $\theta$  is  $(-1.27, 1.93)$ . Calculation gives

$$\begin{aligned} \inf_{\pi \in \Gamma} P^\pi(-1.27 < \theta < 1.93 | X = 0.5) &= 0.817, \\ \sup_{\pi \in \Gamma} P^\pi(-1.27 < \theta < 1.93 | X = 0.5) &= 0.966. \end{aligned}$$

Hence, the standard credible set is reasonably robust. On the other hand, suppose  $x = 4$  is observed. [Note that, since  $m(x|\pi_0)$  is  $\mathcal{N}(0, 3)$ , this is not an "outrageous" observation.] Then the usual 95% credible set is  $(1.07, 4.27)$ . However, in this case we have that

$$\begin{aligned} \inf_{\pi \in \Gamma} P^\pi(1.07 < \theta < 4.27 | X = 4) &= 0.1355, \\ \sup_{\pi \in \Gamma} P^\pi(1.07 < \theta < 4.27 | X = 4) &= 0.99. \end{aligned}$$

Since the posterior probability can get as low as 0.1355 for  $x = 4$ , robustness is not present.

Two interesting general points emerge from the previous example. First, robustness with respect to  $\Gamma$  will usually depend significantly on the  $x$  observed. Second, a lack of robustness may be due to the fact that  $\Gamma$  is "too large." When  $x = 4$ , for instance, the low probability of coverage (0.1355) is achieved when the contamination,  $q$ , is a point mass at 4.27. The resulting prior would probably not have been deemed to be reasonable a priori. Using a more reasonable  $\gamma$  might result in robustness. Also, more robust credible sets can be found—indeed Berger and Berliner (1983) determine the optimal  $1 - \alpha$  robust credible set, optimal in the sense of having smallest size (Lebesgue measure) subject to its posterior probability being at least  $1 - \alpha$  for all  $\pi$  in  $\Gamma$ . In any case, the use of  $\mathcal{Q} = \mathcal{P}$  and Theorem 2.2 is conservative, in that, if robustness of a credible set is achieved for such  $\Gamma$ , one knows that robustness is also present for the more reasonable, smaller  $\Gamma$ .

Theorem 2.2 can also be used for hypothesis testing. Thus suppose we desire to test the hypothesis  $H_0: \theta \in \Theta_0$  versus the alternative  $H_1: \theta \in \Theta - \Theta_0$ . For a fixed prior  $\pi$ , the usual Bayesian test is based on the posterior odds ratio  $O_\pi(x)$ , defined by

$$O_\pi(x) = P^\pi(\theta \in \Theta_0 | X = x) / [1 - P^\pi(\theta \in \Theta_0 | X = x)].$$

Letting  $C = \Theta_0$ , Theorem 2.2 immediately yields the following:

**COROLLARY 2.1.** For  $\Gamma$  as in (2.1),

$$\inf_{\pi \in \Gamma} O_\pi(x) = O_{\pi_0}(x) \left\{ 1 + \frac{\varepsilon \sup_{\theta \notin \Theta_0} f(x|\theta)}{(1 - \varepsilon)(1 - \beta_0)m(x|\pi_0)} \right\}^{-1},$$

and

$$\sup_{\pi \in \Gamma} O_\pi(x) = O_{\pi_0}(x) \left\{ 1 + \frac{\varepsilon \sup_{\theta \in \Theta_0} f(x|\theta)}{(1 - \varepsilon)(1 - \beta_0)m(x|\pi_0)} \right\},$$

where  $\beta_0 = P^{\pi_0}(\theta \in \Theta_0 | X = x)$ .

In testing, it will usually be much easier to achieve robustness using this “too large”  $\Gamma$ , since extreme  $x$  [i.e.,  $x$  for which  $m(x|\pi_0)$  is small], which lead to the unrealistic point mass contaminations, will usually provide extreme evidence for, or against,  $\Theta_0$ . (The difference between the inf and sup of  $O_\pi$  may be substantial, but they will both be substantially less than one or substantially greater than one.) Together with the simplicity of the results in Corollary 2.1, this makes the use of  $\mathcal{L} = \mathcal{P}$  very attractive for robustness investigations in testing.

It should be clear that Theorem 2.2 is also immediately applicable to the testing of several hypotheses and to classification problems. Lower and upper bounds on the posterior probabilities of all hypotheses can be obtained.

**3. Symmetric unimodal contaminations.** A natural, yet remarkably tractable, class of priors to consider when  $\Theta \subseteq \mathbb{R}^1$ , is the  $\epsilon$ -contamination class defined (for fixed  $\theta_0$ ) by

$$(3.1) \quad \mathcal{L} = \{ \text{densities of the form } q(|\theta - \theta_0|), q \text{ nonincreasing} \}.$$

This class is particularly reasonable when  $\pi_0$  itself is unimodal and symmetric about  $\theta_0$ . Note that under such circumstances, the resulting contaminated priors  $\pi$  display the desirable properties that (i) values of  $\theta$  far from  $\theta_0$  cannot be given overwhelming weight (unlike the possibilities observed in Section 2), but (ii) priors with tails larger than  $\pi_0$  are considered.

The considerable simplicity of working with (3.1) accrues from the fact that in much of the analysis, (3.1) can be replaced by

$$(3.2) \quad \mathcal{L}' = \{ \text{Uniform } (\theta_0 - a, \theta_0 + a) \text{ densities, } a \geq 0 \},$$

where the “density” when  $a = 0$  is a point mass at  $\theta_0$ . Required optimizations thus involve only the variable  $a$ . Preliminary analyses of the type discussed in Section 2.4 have been carried out in this manner, but are a bit involved and will be reported elsewhere. However, the ML-II prior is quite simple to present:

**THEOREM 3.1.** *For the  $\epsilon$ -contamination class with  $\mathcal{L}$  as in (3.1), an ML-II prior is*

$$\hat{\pi} = (1 - \epsilon)\pi_0 + \epsilon\hat{q},$$

where  $\hat{q}$  is Uniform  $(\theta_0 - \hat{a}, \theta_0 + \hat{a})$ ,  $\hat{a}$  being the value of  $a$  which maximizes

$$(3.3) \quad m(x|a) = \begin{cases} (2a)^{-1} \int_{\theta_0-a}^{\theta_0+a} f(x|\theta) d\theta, & a > 0, \\ f(x|\theta_0), & a = 0. \end{cases}$$

**PROOF.** The proof follows trivially after noting that (i) any prior in (3.1) is a mixture of priors in (3.2), and (ii)  $m(x|\pi)$  is a linear functional of  $\pi$ .  $\square$

Theorem 3.1 is an adaptation of a result in Berger and Sellke (1984), who utilized the fact that  $m(x|\hat{a})$  is an upper bound on  $m(x|q)$ ,  $q$  in  $\mathcal{L}$  given by (3.1), to establish startling lower bounds on posterior probabilities of point null

hypotheses that are an order of magnitude larger than classical significance levels or  $P$ -values. Here we utilize the theorem to calculate the ML-II posterior mean and variance in an illustrative example:

**EXAMPLE: ESTIMATING A NORMAL MEAN.** Suppose  $X \sim \mathcal{N}(\theta, \sigma^2)$ ,  $\sigma^2$  known, and  $\pi_0$  is  $\mathcal{N}(\theta_0, \tau^2)$ ,  $\theta_0$  and  $\tau^2$  given. Let  $\mathcal{Q}$  be as in (3.1) and define  $z = (x - \theta_0)/\sigma$ . Following Berger and Sellke (1984) [see also Berger (1985)],  $\hat{a}$  of Theorem 3.1 is

$$(3.4) \quad \hat{a} = \begin{cases} 0 & \text{if } z \leq 1.65, \\ \alpha^* \sigma & \text{if } z > 1.65, \end{cases}$$

where  $\alpha^*$  satisfies the equation

$$(3.5) \quad \hat{a}^* = |z| + \left[ -2 \log \left( \sqrt{2\pi} \left\{ [\Phi(\alpha^* - |z|) - \Phi(-(\alpha^* + |z|))]/\alpha^* - \phi(-(\alpha^* + |z|)) \right\} \right) \right]^{1/2},$$

$\Phi$  and  $\phi$  denoting the standard normal c.d.f. and density, respectively. Equation (3.5) can be solved (usually very quickly) by standard fixed point iteration, starting on the right-hand side with initial value  $\alpha^* = |z|$ .

**CASE 1:  $\hat{a} > 0$ .** It can be shown that the ML-II posterior mean and variance corresponding to the Uniform ( $\theta_0 - \hat{a}$ ,  $\theta_0 + \hat{a}$ ) prior are, respectively,

$$\delta^{\hat{q}}(x) = x - (\sigma/\alpha^*) \tanh(z\alpha^*)$$

and

$$V^{\hat{q}}(x) = \sigma^2 \left[ z - (\sigma^2/\alpha^*) \tanh(z\alpha^*) \right] \left[ \alpha^{*-1} \tanh(z\alpha^*) \right].$$

Furthermore,

$$\hat{\lambda}(x) = \left[ 1 + \left( 0.5\epsilon(1 + \tau^2/\sigma^2)^{1/2}/(1 - \epsilon) \right) (1 + \exp(-2z\alpha^*)) \right. \\ \left. \times \exp \left\{ -0.5 \left( (\tau^2 z^2 / (\sigma^2 + \tau^2)) + \alpha^{*2} - 2z\alpha^* \right) \right\} \right]^{-1},$$

so that the ML-II posterior mean and variance are [see (1.7) and (1.8)]

$$(3.6) \quad \delta^{\hat{\pi}}(x) = \hat{\lambda}(x) \left( x - (\sigma^2/(\sigma^2 + \tau^2))(x - \theta_0) \right) + (1 - \hat{\lambda}(x)) \delta^{\hat{q}}(x)$$

and

$$(3.7) \quad V^{\hat{\pi}}(x) = \hat{\lambda}(x) \tau^2 \sigma^2 / (\tau^2 + \sigma^2) + (1 - \hat{\lambda}(x)) V^{\hat{q}}(x) \\ + \hat{\lambda}(x) (1 - \hat{\lambda}(x)) \left[ (\sigma/\alpha^*) \tanh(z\alpha^*) - \sigma^2(x - \theta_0)/(\sigma^2 + \tau^2) \right]^2.$$

**CASE 2:  $\hat{a} = 0$ .** If  $\hat{a} = 0$ , then  $\delta^{\hat{q}}(x) = \theta_0$  and  $V^{\hat{q}}(x) = 0$ . The formulas for  $\hat{\lambda}$ ,  $\delta^{\hat{\pi}}$ , and  $V^{\hat{\pi}}$  are then easily obtained and omitted here. In this case, however, though  $x$  is close to  $\theta_0$ , it is probably undesirable to allow  $\hat{\pi}$  to be more concentrated at  $\theta_0$  than  $\pi_0$ . The natural "fix-up" is to simply replace  $(\delta^{\hat{\pi}}, V^{\hat{\pi}})$  by  $(\delta^{\pi_0}, V^{\pi_0})$ . In fact we generally recommend this modification whenever  $V^{\hat{\pi}} < V^{\pi_0}$ .

TABLE 1  
ML-II results for unimodal, symmetric contaminations

$x$	$\hat{a}$	$\delta^{\pi_0}$	$\hat{\lambda}$	$\delta^{\hat{\pi}}$	$V^{\hat{\pi}}$
3.00	4.13	2.000	0.661	2.257	0.796
5.00	6.43	3.333	0.166	4.594	1.054
7.00	8.61	4.667	$4.7 \times 10^{-3}$	6.873	0.842
10.00	11.78	6.667	$1.3 \times 10^{-6}$	9.915	0.851

As a specific example, suppose  $\sigma^2 = 1$ ,  $\theta_0 = 0$ ,  $\tau^2 = 2$ , and  $\varepsilon = 0.2$ . Values of ML-II quantities and  $\delta^{\pi_0}(x) = 2x/3$  for various  $x$  are given in Table 1. Further numerical results are given in Section 6.

**4. Unimodality preserving contaminations.**

*4.1. Introduction.* Despite the simplicity of the analyses in Sections 2 and 3, there are some objections. In Section 2 we saw that choosing  $\mathcal{Q} = \mathcal{P} = \{\text{all distributions}\}$  can cause serious problems, due to the unrealistic nature of some of the resultant priors in  $\Gamma$ . On the other hand, the restriction to symmetric unimodal contaminations in Section 3 could be criticized for not allowing certain plausible contaminations, particularly when  $\pi_0$  is not symmetric. When  $\Theta \subset \mathbb{R}^1$  and  $\pi_0$  is unimodal, as will be assumed throughout this section, perhaps *the most appealing*  $\mathcal{Q}$  is that which contains those contaminations which preserve the unimodality of  $\pi = (1 - \varepsilon)\pi_0 + \varepsilon q$  (note that  $q$  need not be unimodal). Any such  $\pi$  would typically be plausible a priori, and virtually *all*  $\pi$  deemed reasonable a priori are in this class. Successfully working with such “minimally complete”  $\Gamma$  is a very desirable goal in Bayesian robustness.

The main goals of this section are to indicate that such complex classes can, surprisingly, be analyzed and to provide some basic mathematical techniques likely to be encountered in such analyses. The presentation here is restricted to the determination of the ML-II prior and the ML-II posterior mean and variance, a detailed example concerning the mean of a normal distribution, plus some discussion of results later in Section 6. We are currently investigating the more fundamental question of determining ranges of posterior measures as  $\pi$  varies over  $\Gamma$ .

The exact class  $\Gamma$  that will be considered is (where  $\theta_0$  denotes the mode of  $\pi_0$ , which we assume to be unique)

$$\Gamma = \{ \pi = (1 - \varepsilon)\pi_0 + \varepsilon q : q \in \mathcal{Q}, \text{ the set of all probability densities for which } \pi \text{ is unimodal with (not necessarily unique) mode } \theta_0, \text{ and } \pi(\theta_0) \leq (1 + \varepsilon')\pi_0(\theta_0) \}.$$

The final inequality in the definition of  $\Gamma$  specifies the reasonable constraint that  $q$  not be allowed to concentrate too sharply near  $\theta_0$ . (Usually it would be reasonable to select  $\varepsilon' = \varepsilon$ , but this is not necessary. Indeed the choice  $\varepsilon' = 0$

might sometimes be desired, it having the attractive property of ensuring that the ML-II posterior variance never drops much below that of  $\pi_0$  because of excessive prior concentration near  $\theta_0$ .) We will also assume that the likelihood function  $f(x|\theta)$  is unimodal (as a function of  $\theta$ , of course) with unique mode  $\hat{\theta}$ . [Of course,  $x$  is fixed, so  $f(x|\theta)$  need only be unimodal for the observed  $x$ , not for all  $x$ .] It will also be technically convenient to restrict consideration to  $\pi_0$  and  $f$  which are nonzero and are strictly monotonic on each side of the modes. More general cases could be handled, but the results get messier. We also assume, without loss of generality, that  $\hat{\theta} \geq \theta_0$ .

4.2. *The form of the ML-II prior and posterior.* The formal calculation, in Sections 4.4 and 4.5, of the ML-II prior,  $\hat{\pi}$ , is complicated by the need to consider several different cases. The result, however, is always of the quite simple form

$$(4.1) \quad \hat{\pi}(\theta) = \begin{cases} K & \text{for } \theta \in B, \\ (1 - \varepsilon)\pi_0(\theta) & \text{for } \theta \notin B, \end{cases}$$

i.e.,  $\hat{\pi}$  is uniform over  $B$  (which will be an interval about  $\hat{\theta}$ ), and equals  $(1 - \varepsilon)\pi_0(\theta)$  outside of  $B$ . (Note that  $K$  is implicitly defined by the constraint that  $\hat{\pi}$  have mass one.) Thus the ML-II posterior will be

$$(4.2) \quad \hat{\pi}(\theta|x) = \hat{\lambda}(x)\pi_0(\theta|x) + (1 - \hat{\lambda}(x))\hat{q}(\theta|x),$$

where [letting  $I_B(\theta)$  denote the usual indicator function on  $B$ ]

$$(4.3) \quad \begin{aligned} \hat{q} &= \varepsilon^{-1} [K - (1 - \varepsilon)\pi_0(\theta)] I_B(\theta), \\ \hat{\lambda}(x) &= (1 - \varepsilon)m(x|\pi_0) / [(1 - \varepsilon)m(x|\pi_0) + \varepsilon m(x|\hat{q})], \end{aligned}$$

$$m(x|\hat{q}) = \varepsilon^{-1} \int_B [K - (1 - \varepsilon)\pi_0(\theta)] f(x|\theta) d\theta.$$

The interesting case (Case 1 in Section 4.5) is that in which  $|\hat{\theta} - \theta_0|$  is moderately large (i.e., where prior and likelihood are not in close agreement), since it is only in this case that the choice of  $\pi \in \Gamma$  will have a substantial effect. As  $|\hat{\theta} - \theta_0| \rightarrow \infty$ , it can typically be shown that the uniform piece of  $\hat{\pi}$  dominates, in the sense that

$$\hat{\pi}(\theta|x) \rightarrow f(x|\theta) / \int f(x|\theta) d\theta,$$

which would be the posterior for the noninformative uniform prior. This kind of behavior can be labelled "robust" from a number of viewpoints [cf. Berger (1984)], and is certainly more pleasing than the limiting behavior of  $\hat{\pi}(\theta|x)$  in Section 2, which collapsed to a point at  $\hat{\theta}$  in the limit. Discussion of the degree of "robustness" which is attained by  $\hat{\pi}$  is delayed until Section 6.

4.3. *The normal distribution.* We present here an example of the overall theory, using the formulas from Sections 4.4 and 4.5. The example considered is that in which  $X$  is  $\mathcal{N}(\theta, 1)$  and  $\pi_0$  is  $\mathcal{N}(0, \tau^2)$ . [The more general case where  $X \sim \mathcal{N}(\theta, \sigma^2)$  and  $\pi_0$  is  $\mathcal{N}(\mu, \tau^2)$ ,  $\sigma^2$ ,  $\tau^2$ , and  $\mu$  all known, can be reduced to

this case by a linear transformation.] Only Case 1 will be considered (which here means that  $|x|$  is larger than a certain constant depending on  $\varepsilon$  and  $\tau$ ) and we assume that  $x > 0$ . Let  $\Phi$  denote the standard normal c.d.f., and  $\phi$  the standard normal density function.

The set  $B$  in (4.1) is the interval  $[\theta^*, w(\theta^*)]$ , the endpoints being implicitly defined by the equations [noting  $\theta^* < x < w(\theta^*)$ ]

$$(4.4) \quad (1 - \varepsilon)\phi\left(\frac{\theta^*}{\tau}\right)\left[\frac{w(\theta^*)}{\tau} - \frac{\theta^*}{\tau}\right] - (1 - \varepsilon)\left[\Phi\left(\frac{w(\theta^*)}{\tau}\right) - \Phi\left(\frac{\theta^*}{\tau}\right)\right] = \varepsilon,$$

$$\phi(w(\theta^*) - x)[w(\theta^*) - \theta^*] - [\Phi(w(\theta^*) - x) - \Phi(\theta^* - x)] = 0.$$

These equations can be easily solved by iteration. Calculation then gives that the ML-II posterior is

$$\hat{\pi}(\theta|x) = \begin{cases} \hat{\lambda}(x)C_0f(x|\theta) & \text{if } \theta \in B, \\ \hat{\lambda}(x)\pi_0(\theta|x) & \text{if } \theta \notin B, \end{cases}$$

where

$$\pi_0(\theta|x) \text{ is } \mathcal{N}(\delta, V^2), \quad \delta = \frac{\tau^2x}{1 + \tau^2}, \quad V^2 = \frac{\tau^2}{1 + \tau^2}$$

(this notation is more convenient here than the previous  $\delta^{\pi_0}, V^{\pi_0}$ ),

$$\hat{\lambda}(x) = [1 - B_0 + C_0\phi(w(\theta^*) - x)(w(\theta^*) - \theta^*)]^{-1},$$

$$C_0 = \frac{(1 + \tau^{-2})^{1/2}\phi(\theta^*/\tau)}{\phi(x[1 + \tau^2]^{-1/2})}, \quad \text{and} \quad B_0 = \Phi\left(\frac{w(\theta^*) - \delta}{V}\right) - \Phi\left(\frac{\theta^* - \delta}{V}\right).$$

The ML-II posterior mean and variance are given, respectively, by

$$\delta^{\hat{\pi}} = \hat{\lambda}(x)\delta + [1 - \hat{\lambda}(x)]\delta^{\hat{q}}$$

and

$$V^{\hat{\pi}} = \hat{\lambda}(x)V^2 + [1 - \hat{\lambda}(x)]V^{\hat{q}} + \hat{\lambda}(x)[1 - \hat{\lambda}(x)][\delta - \delta^{\hat{q}}]^2,$$

where

$$\delta^{\hat{q}} = x + (C_0D_0 - E_0 + (x - \delta)B_0)\frac{\hat{\lambda}}{(1 - \hat{\lambda})},$$

$$V^{\hat{q}} = [\hat{\lambda}(x)^{-1} - 1]^{-1}\left\{-C_0D_0(2\delta^{\hat{q}} - x - \theta^*) + (x - \delta^{\hat{q}})^2[\hat{\lambda}(x)^{-1} - 1 + B_0] - (\theta^* + \delta - 2\delta^{\hat{q}})E_0 - [V^2 + (\delta - \delta^{\hat{q}})^2]B_0 + V[w(\theta^*) - \theta^*]\phi([w(\theta^*) - \delta]/V)\right\},$$

$$D_0 = \phi(\theta^* - x) - \phi(w(\theta^*) - x),$$

and

$$E_0 = V\left\{\phi\left(\frac{\theta^* - \delta}{V}\right) - \phi\left(\frac{w(\theta^*) - \delta}{V}\right)\right\}.$$

TABLE 2  
ML-II quantities for various  $x$

$x$	$B$	$\delta$	$\hat{\lambda}$	$\delta^\dagger$	$V^\dagger$
1.75	(0, 2.53)	1.167	0.609	1.425	0.599
3.00	(1.39, 3.73)	2.000	0.375	2.581	0.616
5.00	(2.25, 6.08)	3.333	0.052	4.735	0.666
7.00	(2.64, 8.35)	4.667	$1.5 \times 10^{-3}$	6.827	0.735
10.00	(2.97, 11.61)	6.667	$4.5 \times 10^{-7}$	9.880	0.797

As a specific example, suppose  $\tau^2 = 2$  and  $\epsilon = 0.2$ . Then one can show [using (4.6)] that Case 1 occurs providing  $|x| \geq 1.75$ . Table 2 gives the relevant quantities above for various  $x$ .

The behavior alluded to earlier clearly obtains: as  $|x|$  gets large,  $\hat{\lambda} \rightarrow 0$ , and the uniform part of the prior (on  $B$ ) dominates. Also,  $\delta^{\hat{\tau}}(x) \rightarrow x$  and  $V^{\hat{\tau}}(x) \rightarrow 1$ . Indeed, the following theorem gives large  $x$  approximations to the key quantities, approximations which are accurate, in the above example, for  $x \geq 7$  and which show that the domination of the uniform portion occurs at an exponential rate. (The proof of the theorem is routine and will be omitted.)

**THEOREM 4.1.** As  $x \rightarrow \infty$  ( $\log$  denotes natural logarithm),

$$\begin{aligned} \theta^* &= [2\tau^2 \log x]^{1/2} + o(1), & w(\theta^*) &= x + [2 \log(x/\sqrt{2\pi})]^{1/2} + o(1), \\ \hat{\lambda}(x) &= (\epsilon^{-1} - 1)x[(1 + \tau^2)2\pi]^{-1/2} \exp\{-x^2/[2(1 + \tau^2)]\}(1 + o(1)), \\ \delta^{\hat{q}}(x) &= x - x^{-1}(1 + o(1)), & V^{\hat{q}}(x) &= 1 - x^{-1}[2 \log x]^{1/2}(1 + o(1)). \end{aligned}$$

**4.4. Preliminaries and notation for the general theory.** For  $-\epsilon' \leq \rho \leq \epsilon$ , define  $v(\rho) \geq \theta_0$ , implicitly, by

$$(4.5) \quad \pi_0(\theta_0)(1 - \rho)(v(\rho) - \theta_0) - (1 - \epsilon) \int_{\theta_0}^{v(\rho)} \pi_0(\theta) d\theta = \epsilon,$$

and define

$$(4.6) \quad V(\rho) = f(x|v(\rho))(v(\rho) - \theta_0) - \int_{\theta_0}^{v(\rho)} f(x|\theta) d\theta.$$

For  $\theta_0 \leq \theta$ , define  $w(\theta) \geq \theta$ , implicitly, by

$$(4.7) \quad (1 - \epsilon)\pi_0(\theta)(w(\theta) - \theta) - (1 - \epsilon) \int_{\theta}^{w(\theta)} \pi_0(\xi) d\xi = \epsilon,$$

and define

$$(4.8) \quad W(\theta) = f(x|w(\theta))(w(\theta) - \theta) - \int_{\theta}^{w(\theta)} f(x|\xi) d\xi.$$

**LEMMA 4.1.** (a) *The quantities  $v(\rho)$  and  $w(\theta)$  are well defined, unique, continuous, and strictly increasing for  $-\epsilon' \leq \rho \leq \epsilon$  and  $\theta \geq \theta_0$ .*

(b) If  $v(\rho) > \hat{\theta}$ , then  $V(\rho)$  is decreasing in  $\rho$ . Furthermore,  $V(\rho) = 0$  has at most one solution.

(c) If  $\theta_0 \leq \theta \leq \hat{\theta}$  and  $w(\theta) > \hat{\theta}$ , then  $W(\theta)$  is decreasing at  $\theta$ . Furthermore, if  $V(\epsilon) \geq 0$ , then  $W(\theta) = 0$  has a unique solution  $\theta_0 \leq \theta^* < \hat{\theta}$ .

PROOF. (a) At  $v = \theta_0$ , the left-hand side of (4.5) is zero. As  $v \rightarrow \infty$ , the left-hand side of (4.5) goes to  $\infty$ . Finally, since  $\pi_0$  is decreasing for  $\theta > \theta_0$ , the derivative, with respect to  $v$ , of the left-hand side of (4.5) is easily seen to be strictly positive. A solution to (4.5) thus exists and is unique.

To show that  $v(\rho)$  is strictly increasing, one can differentiate both sides of (4.5) with respect to  $\rho$  and solve for  $v'(\rho)$  [i.e.,  $d/d\rho v(\rho)$ ], obtaining

$$v'(\rho) = \pi_0(\theta_0)(v(\rho) - \theta_0) / [\pi_0(\theta_0)(1 - \rho) - (1 - \epsilon)\pi_0(v(\rho))].$$

Since  $v(\rho) > \theta_0$ ,  $\rho < \epsilon$ , and  $\pi_0$  is decreasing for  $\theta > \theta_0$ , it is clear that  $v'(\rho) > 0$ . The verification for  $w(\theta)$  is very similar.

(b) Letting  $f'(x|\theta) = d/d\theta f(x|\theta)$ , calculation gives

$$\frac{d}{d\theta} V(\rho) = f'(x|v(\rho))v'(\rho)(v(\rho) - \theta_0).$$

Since  $f$  is decreasing for  $\theta > \hat{\theta}$ , the monotonicity result follows from part (a). If  $V(\rho) = 0$ , the unimodality of  $f$  ensures that  $v(\rho) > \hat{\theta}$  [for otherwise the right-hand side of (4.6) is positive]. The strict monotonicity of  $V$  for such  $\rho$  ensures that any solution to  $V(\rho) = 0$  must be unique.

(c) Letting  $w'(\theta) = d/d\theta w(\theta)$ , calculation gives

$$\frac{d}{d\theta} W(\theta) = f'(x|w(\theta))w'(\theta)(w(\theta) - \theta).$$

The monotonicity of  $f$  and part (a) show that this is negative. Using this, to show that  $W(\theta) = 0$  has a unique solution, it is only necessary to show that  $W(\theta_0) \geq 0$  and  $W(\hat{\theta}) < 0$ . Since  $v(\epsilon) = w(\theta_0)$ , it follows that  $W(\theta_0) = V(\epsilon) \geq 0$  (by assumption). That  $W(\hat{\theta}) < 0$  follows from (4.8) and an easy application of the mean value theorem [since  $f(x|\theta)$  decreases for  $\theta > \hat{\theta}$ ].  $\square$

LEMMA 4.2. Suppose  $V(\epsilon) \geq 0$ , and let  $\theta_0 \leq \theta^* \leq \hat{\theta}$  be the solution to  $W(\theta) = 0$ . Then

(a)  $f(x|\theta) < f(x|w(\theta^*))$  for  $\theta \notin [\theta^*, w(\theta^*)]$ .

(b) For any nonincreasing integrable function  $g$  such that  $\int g(\theta) d\theta = 0$ , it follows that

$$(4.9) \quad \int_{\theta^*}^{w(\theta^*)} g(\theta) f(x|\theta) d\theta \leq 0.$$

PROOF. (a) Clearly  $f(x|\theta^*) < f(x|w(\theta^*))$ , for otherwise the integrand in (4.8) would be everywhere larger than  $f(x|w(\theta^*))$  and  $W(\theta^*)$  would be nonzero, a contradiction. The unimodality of  $f$  thus gives the result for  $\theta < \theta^*$ . Now  $w(\theta^*) > \hat{\theta}$ , for otherwise (4.8) could again be used to contradict  $W(\theta^*) = 0$ . The unimodality of  $f$  thus also gives the result for  $\theta > w(\theta^*)$ .



(b) Note first that it suffices to prove the result for differentiable  $g$ . Letting  $h(\theta) = -d/d\theta g(\theta)$  (note  $h \geq 0$ ) and writing  $g(\theta) = K - \int_{\theta^*}^{\theta} h(\xi) d\xi$ , where

$$(4.10) \quad \begin{aligned} K &= \frac{1}{[w(\theta^*) - \theta^*]} \int_{\theta^*}^{w(\theta^*)} \int_{\theta^*}^{\eta} h(\xi) d\xi d\eta \\ &= \frac{1}{[w(\theta^*) - \theta^*]} \int_{\theta^*}^{w(\theta^*)} (w(\theta^*) - \xi)h(\xi) d\xi, \end{aligned}$$

we obtain from Fubini's theorem

$$(4.11) \quad \begin{aligned} \int_{\theta^*}^{w(\theta^*)} g(\theta) f(x|\theta) d\theta &= K \int_{\theta^*}^{w(\theta^*)} f(x|\theta) d\theta \\ &\quad - \int_{\theta^*}^{w(\theta^*)} h(\xi) \int_{\xi}^{w(\theta^*)} f(x|\theta) d\theta d\xi. \end{aligned}$$

Next we show that, for  $\theta^* < \xi < w(\theta^*)$ ,

$$(4.12) \quad \psi(\xi) \equiv \int_{\xi}^{w(\theta^*)} f(x|\theta) d\theta \geq (w(\theta^*) - \xi) f(x|w(\theta^*)).$$

For  $\xi \geq \hat{\theta}$  this is a trivial consequence of the monotonicity of  $f$ . For  $\theta^* < \xi < \hat{\theta}$ , note that  $\psi(\xi)$  is concave [ $f(x|\xi)$  is increasing here] and that

$$(4.13) \quad \psi(\theta^*) = \int_{\theta^*}^{w(\theta^*)} f(x|\theta) d\theta = (w(\theta^*) - \theta^*) f(x|w(\theta^*))$$

[since  $W(\theta^*) = 0$ ]. Hence,  $\psi(\xi)$  must lie above the line  $(w(\theta^*) - \xi) f(x|w(\theta^*))$ , establishing (4.12). Using (4.12) in (4.11) we get that

$$\begin{aligned} &\int_{\theta^*}^{w(\theta^*)} g(\theta) f(x|\theta) d\theta \\ &\leq K \int_{\theta^*}^{w(\theta^*)} f(x|\theta) d\theta - f(x|w(\theta^*)) \int_{\theta^*}^{w(\theta^*)} (w(\theta^*) - \xi)h(\xi) d\xi, \end{aligned}$$

the right-hand side of which is zero by (4.10) and (4.13).  $\square$

4.5. *The ML-II prior.* Define  $\hat{\pi}$  as follows:

CASE 1. If  $V(\varepsilon) \geq 0$ , and  $\theta^* \in [\theta_0, \hat{\theta}]$  is the solution to  $W(\theta) = 0$ , let

$$(4.14) \quad \hat{\pi}(\theta) = \begin{cases} (1 - \varepsilon)\pi_0(\theta^*) & \text{for } \theta^* \leq \theta \leq w(\theta^*), \\ (1 - \varepsilon)\pi_0(\theta) & \text{otherwise.} \end{cases}$$

CASE 2. If  $V(\varepsilon) < 0$  but  $V(-\varepsilon') \geq 0$ , find  $\rho^* \in [-\varepsilon', \varepsilon]$  so that  $V(\rho^*) = 0$ , and let

$$(4.15) \quad \hat{\pi}(\theta) = \begin{cases} (1 - \rho^*)\pi_0(\theta_0) & \text{for } \theta_0 \leq \theta \leq v(\rho^*), \\ (1 - \varepsilon)\pi_0(\theta) & \text{otherwise.} \end{cases}$$

CASE 3. If  $V(-\varepsilon') < 0$  and  $f(x|\theta_0) \leq f(x|v(-\varepsilon'))$ , let  $\hat{\pi}$  be as in Case 2 with  $\rho^* = -\varepsilon'$ .

CASE 4. If  $V(-\epsilon') < 0$  and  $f(x|\theta_0) > f(x|v(-\epsilon'))$ , let

$$(4.16) \quad \hat{\pi}(\theta) = \begin{cases} (1 + \epsilon')\pi_0(\theta_0) & \text{for } \theta' \leq \theta \leq \theta'', \\ (1 - \epsilon)\pi_0(\theta) & \text{otherwise,} \end{cases}$$

where  $\theta'$  and  $\theta''$  are the (unique) solutions to the equations

$$(4.17) \quad \begin{aligned} f(x|\theta') &= f(x|\theta''), \\ (1 + \epsilon')\pi_0(\theta_0)(\theta'' - \theta') - (1 - \epsilon) \int_{\theta'}^{\theta''} \pi(\theta) d\theta &= \epsilon. \end{aligned}$$

Lemma 4.1 establishes that all quantities involved in the definition of  $\hat{\pi}$  are well defined and unique. (The existence and uniqueness of  $\theta'$  and  $\theta''$  in Case 4 is easy to establish.) Observe that, in all cases,  $\hat{\pi}$  has a very simple and easy to work with form of being uniform in a certain interval, and otherwise being equal to  $(1 - \epsilon)\pi_0$ . Case 1 corresponds to the situation where the elicited prior,  $\pi_0$ , and the likelihood function,  $f(x|\theta)$ , are moderately separated, Case 2 to the situation where they are fairly close, and Cases 3 and 4 to situations where they are very close.

**THEOREM 4.2.** *The  $\hat{\pi}$  defined in (4.14)–(4.16) is the ML-II prior in  $\Gamma$ .*

**PROOF.** We only present the argument for Case 1, the other cases being very similar. The goal is to show that

$$(4.18) \quad m(x|\pi) - m(x|\hat{\pi}) = \int [\pi(\theta) - \hat{\pi}(\theta)] f(x|\theta) d\theta \leq 0$$

for all  $\pi \in \Gamma$ . Letting  $g(\theta) = \pi(\theta) - \hat{\pi}(\theta)$ , note that

- (i)  $g(\theta) \geq 0$  for  $\theta \notin [\theta^*, w(\theta^*)]$ , since  $\hat{\pi}(\theta) = (1 - \epsilon)\pi_0(\theta)$  here and  $\pi(\theta) \geq (1 - \epsilon)\pi_0(\theta)$ ;
- (ii)  $g(\theta)$  is nonincreasing on  $[\theta^*, w(\theta^*)]$ , since  $\hat{\pi}(\theta)$  is uniform on this interval and so  $\pi(\theta) = g(\theta) + \hat{\pi}(\theta)$  would have a secondary mode were  $g(\theta)$  somewhere increasing;
- (iii)  $K \equiv \int_{\theta^*}^{w(\theta^*)} g(\theta) d\theta = - \int_{[\theta^*, w(\theta^*)]^c} g(\theta) d\theta$ .

Lemma 4.2(a) and (i) show that

$$\int_{[\theta^*, w(\theta^*)]^c} g(\theta) f(x|\theta) d\theta < f(x|w(\theta^*))(-K).$$

Lemma 4.2(b) and (ii) imply that

$$\int_{\theta^*}^{w(\theta^*)} \left( g(\theta) - \frac{K}{[w(\theta^*) - \theta^*]} \right) f(x|\theta) d\theta \leq 0.$$

Thus

$$(4.19) \quad \int g(\theta) f(x|\theta) d\theta < f(x|w(\theta^*))(-K) + \frac{K}{[w(\theta^*) - \theta^*]} \int_{\theta^*}^{w(\theta^*)} f(x|\theta) d\theta.$$

Since  $W(\theta^*) = 0$ , the right-hand side of (4.19) is zero, and (4.18) follows.  $\square$

COMMENTS. 1. The key step in the proof of Theorem 4.2 is really Lemma 4.2(b), which shows that one cannot improve on a uniform  $\hat{\pi}$  on  $[\theta^*, \omega(\theta^*)]$ . 2. The problem might be susceptible to attack through calculus of variations, since one is trying to maximize an expression involving an integral of  $\pi$  over a class of  $\pi$ . The difficulty is that the  $\pi \in \Gamma$  satisfy a large number of inequality and differential inequality constraints. Calculus of variations with such side constraints is quite difficult.

## 5. Hierarchical classes of priors.

5.1. *Introduction.* Hierarchical priors are typically employed when  $\theta$  is a vector  $(\theta_1, \theta_2, \dots, \theta_p)$ , and the  $\theta_i$  are thought to be independent realizations from a common prior distribution  $g$ . Typically  $g$  is assumed to lie in some class  $\Gamma_1 = \{g_\omega: \omega \in \Omega\}$  of distributions, often the class of conjugate priors, and a "second stage" prior  $h_0$  is placed on this class, i.e., on  $\omega$ . Such a hierarchical prior can, of course, be written as a single prior, namely

$$(5.1) \quad \pi_0(\theta) = \int_{\Omega} \left[ \prod_{i=1}^p g_{\omega}(\theta_i) \right] h_0(\omega) d\omega.$$

(We restrict ourselves to densities in this section, for convenience, and also will not consider hierarchical priors with more than two stages.) Development of and references for this approach can be found in Good (1980), Lindley and Smith (1972), Morris (1983), and Berger (1985).

There are three possible robustness concerns in working with (5.1). One could question the assumptions (i) that the  $\theta_i$  are i.i.d.; (ii) that the prior  $g$  belongs to  $\Gamma_1$ ; and (iii) that  $h_0$  is specified correctly. Each of these concerns deserves careful consideration separately but in the following we will simply deal with uncertainty in the second stage (i.e.,  $h_0$ ), or in both the first and second stages together.

Simultaneous uncertainty in different stages or aspects of a prior can often be expressed most simply by allowing more than one contamination in the  $\varepsilon$ -contamination model. For instance, one could consider

$$(5.2) \quad \Gamma = \{ \pi: \pi = (1 - \varepsilon_1 - \varepsilon_2)\pi_0 + \varepsilon_1 q_1 + \varepsilon_2 q_2, q_1 \in \mathcal{Q}_1, q_2 \in \mathcal{Q}_2 \},$$

where  $\mathcal{Q}_1$  and  $\mathcal{Q}_2$  are appropriate possible classes of contaminations. Such an extension of the  $\varepsilon$ -contamination model vastly increases its flexibility while causing no real hardship in many applications, because the important formulas (1.4), (1.7), and (1.8) become simply

$$(5.3) \quad m(x|\pi) = (1 - \varepsilon_1 - \varepsilon_2)m(x|\pi_0) + \varepsilon_1 m(x|q_1) + \varepsilon_2 m(x|q_2),$$

$$(5.4) \quad \delta^\pi(x) = [1 - \lambda_1(x) - \lambda_2(x)]\delta^{\pi_0}(x) + \lambda_1(x)\delta^{q_1}(x) + \lambda_2(x)\delta^{q_2}(x),$$

$$(5.5) \quad V^\pi(x) = (1 - \lambda_1 - \lambda_2)V^{\pi_0} + \lambda_1 V^{q_1} + \lambda_2 V^{q_2} + \lambda_1 \lambda_2 (\delta^{q_1} - \delta^{q_2})^2 \\ + (1 - \lambda_1 - \lambda_2)\lambda_1 (\delta^{\pi_0} - \delta^{q_1})^2 + (1 - \lambda_1 - \lambda_2)\lambda_2 (\delta^{\pi_0} - \delta^{q_2})^2,$$

where  $\lambda_i(x) = \varepsilon_i m(x|q^i)/m(x|\pi)$  for  $i = 1, 2$ . Thus one can find the ML-II prior

by separately maximizing  $m(x|q_1)$  and  $m(x|q_2)$  in (5.3) (unless  $\mathcal{Q}_1$  and  $\mathcal{Q}_2$  are related in some fashion) and then easily calculate the resultant ML-II posterior mean and variance.

Before proceeding, it is worthwhile to note that  $\Gamma$  of the form (5.2) might be of interest in other than hierarchical prior situations. Indeed, whenever one has several possible models in mind for the contamination, or even for  $\pi_0$  itself, the uncertainty can be reasonably represented by such a  $\Gamma$ .

*5.2. Second stage uncertainty.* Suppose, in the situation of Section 5.1, that only  $h_0$  is deemed uncertain. (Knowledge at higher levels of hierarchical priors will often be more vague than at lower levels.) An  $\varepsilon$ -contamination model for  $h$  would be

$$(5.6) \quad h(\omega) = (1 - \varepsilon)h_0(\omega) + \varepsilon s(\omega), \quad s \in \mathcal{S}.$$

The resulting prior for  $\theta$  is

$$(5.7) \quad \pi(\theta) = \int \left[ \prod_{i=1}^p g_{\omega}(\theta_i) \right] h(\omega) d\omega = (1 - \varepsilon)\pi_0(\theta) + \varepsilon q(\theta),$$

where

$$\pi_0(\theta) = \int \left[ \prod_{i=1}^p g_{\omega}(\theta_i) \right] h_0(\omega) d\omega \quad \text{and} \quad q(\theta) = \int \left[ \prod_{i=1}^p g_{\omega}(\theta_i) \right] s(\omega) d\omega.$$

Letting  $\mathcal{Q} = \{q: s \in \mathcal{S}\}$ , it follows that the uncertainty in  $\pi$  can be expressed by  $\Gamma = \{\pi: \pi = (1 - \varepsilon)\pi_0 + \varepsilon q, q \in \mathcal{Q}\}$ .

In determining the ML-II prior for this situation, it will be convenient to define

$$m(x|\omega) = \int f(x|\theta) \left[ \prod_{i=1}^p g_{\omega}(\theta_i) \right] d\theta,$$

which is clearly the marginal distribution of  $X$  under the assumption that the prior for  $\theta$  is  $[\prod_{i=1}^p g_{\omega}(\theta_i)]$ . Note that

$$(5.8) \quad m(x|\pi) = (1 - \varepsilon)m(x|\pi_0) + \varepsilon \int m(x|\omega) s(\omega) d\omega.$$

When  $\mathcal{S} = \mathcal{P} = \{\text{all distributions}\}$ , it is clear from (5.8) that

$$\sup_{\pi \in \Gamma_1} m(x|\pi) = (1 - \varepsilon)m(x|\pi_0) + \varepsilon \sup_{\omega} m(x|\omega).$$

Assuming that  $m(x|\omega)$  has a maximum at  $\hat{\omega}$ , it follows that the ML-II prior is

$$\hat{\pi}(\theta) = (1 - \varepsilon)\pi_0(\theta) + \varepsilon \left[ \prod_{i=1}^p g_{\hat{\omega}}(\theta_i) \right],$$

for which analysis is usually quite straightforward.

**EXAMPLE 3.** Suppose that  $X = (X_1, \dots, X_p) \sim \mathcal{N}_p(\theta, \sigma^2 I_p)$ ,  $\sigma^2$  known, and that the first-stage prior information is that the  $\theta_i$  are independent with a

common  $\mathcal{N}(\mu, \tau^2)$  distribution, to be denoted  $g_\omega$ , with  $\omega = (\mu, \tau^2)$  unknown. Note that  $m(x|\omega)$  is  $\mathcal{N}_p(\mu\mathbf{1}, (\sigma^2 + \tau^2)I_p)$ , where  $\mathbf{1} = (1, \dots, 1)$ . It is easy to check that  $m(x|\omega)$  is maximized at

$$\hat{\omega} = (\hat{\mu}, \hat{\tau}^2) = \left( \bar{x}, \max \left[ 0, \frac{1}{p} \sum_{i=1}^p (x_i - \bar{x})^2 - \sigma^2 \right] \right).$$

Hence, with contaminated second-stage prior as in (5.6) and  $\mathcal{S} = \mathcal{P}$ , the ML-II prior is  $\hat{\pi}(\theta) = (1 - \varepsilon)\pi_0(\theta) + \varepsilon\hat{q}(\theta)$ , where  $\hat{q}$  is  $\mathcal{N}_p(\hat{\mu}\mathbf{1}, \hat{\tau}^2 I_p)$ .

As a very special case, suppose  $h_0$  is a point mass at  $(\mu_0, \tau_0^2)$ , so that  $\pi_0$  is simply  $\mathcal{N}_p(\mu_0\mathbf{1}, \tau_0^2 I_p)$ . Then the ML-II posterior is

$$\hat{\pi}(\theta|x) = \hat{\lambda}(x)\pi_0(\theta|x) + (1 - \hat{\lambda}(x))\hat{q}(\theta|x),$$

where  $\pi_0(\theta|x)$  is  $\mathcal{N}_p(\delta^{\pi_0}(x), v_0 I_p)$ ,  $\hat{q}(\theta|x)$  is  $\mathcal{N}_p(\delta^{\hat{q}}(x), \hat{v} I_p)$ ,  $v_0 = \sigma^2 \tau_0^2 / (\sigma^2 + \tau_0^2)$ ,  $\hat{v} = \sigma^2 \hat{\tau}^2 / (\sigma^2 + \hat{\tau}^2)$ ,

$$\delta^{\pi_0}(x) = x - \frac{\sigma^2}{\sigma^2 + \tau_0^2}(x - \mu_0\mathbf{1}), \quad \delta^{\hat{q}}(x) = x - \frac{\sigma^2}{\sigma^2 + \hat{\tau}^2}(x - \hat{\mu}\mathbf{1}),$$

and

$$\begin{aligned} \hat{\lambda}(x) &= (1 - \varepsilon)m(x|\pi_0) / [(1 - \varepsilon)m(x|\pi_0) + \varepsilon m(x|\hat{q})] \\ &= \left[ 1 + \frac{\varepsilon}{(1 - \varepsilon)} \cdot (\sigma^2 + \tau_0^2)^{p/2} \exp \left\{ \sum_{i=1}^p (x_i - \mu_0)^2 / [2(\sigma^2 + \tau_0^2)] \right\} \rho(x) \right]^{-1}, \end{aligned}$$

where

$$\rho(x) = \begin{cases} \sigma^{-p} \exp \left\{ - \sum_{i=1}^p (x_i - \bar{x})^2 / 2p \right\} & \text{if } \sum (x_i - \bar{x})^2 < p\sigma^2, \\ \left[ \frac{1}{p} \sum_{i=1}^p (x_i - \bar{x})^2 \right]^{-p/2} \exp \left\{ - \frac{1}{2}p \right\} & \text{otherwise.} \end{cases}$$

Note that  $\delta^{\pi_0}$  is the usual conjugate prior estimate of  $\theta$ , while  $\sigma^{\hat{q}}$  is the usual empirical Bayes estimate of  $\theta$ . The overall posterior mean [see (1.7)] is thus

$$\delta^{\hat{\pi}} = \hat{\lambda}(x)\delta^{\pi_0}(x) + (1 - \hat{\lambda}(x))\delta^{\hat{q}}(x),$$

which will be close to  $\delta^{\pi_0}$  if the  $x_i$  are close to  $\mu_0$ , and close to  $\delta^{\hat{q}}$  if the  $x_i$  are similar but far from  $\mu_0$ .

Of course, only rarely will it be appropriate to choose  $h_0$  to be a point mass. More natural would be a choice such as  $h_0(\mu, \tau^2) = w(\mu)v(\tau^2)$ , where  $w(\mu)$  is  $\mathcal{N}(\mu_0, A)$  and  $v$  is, say, a gamma distribution. Although the ML-II posterior, is no longer expressible in closed form for such a situation, the posterior mean and variance can be written in a form involving a single numerical integral over  $\tau^2$  [see, e.g., Lindley (1971)].

Several features of the above example are worth noting. First, the strong relationship of the ML-II theory with standard empirical Bayes analysis is apparent. Indeed, if one were to choose  $\varepsilon = 1$ , the standard empirical Bayes

situation would result. As mentioned in the introduction, we much prefer the analysis with reasonably small  $\epsilon$ , the choice  $\epsilon = 1$  resulting (typically) in there being a large number of unrealistic priors in  $\Gamma$ . Of course, the choice  $\mathcal{L} = \mathcal{P}$  also suffers somewhat from this deficiency, as discussed in Section 2.3. An appealing possibility in the above example is, therefore, to attempt to apply the ideas of Section 3 (or possibly Section 4) and work with more reasonable  $\mathcal{L}$ . For instance, if independence of  $\mu$  and  $\tau^2$  can be assumed, so that  $h(\mu, \tau^2) = w(\mu)v(\tau^2)$ , one could elicit  $w_0$  and  $v_0$ , consider

$$\begin{aligned} \mathcal{W} &= \{w = (1 - \epsilon_1)w_0 + \epsilon_1q_w: q_w \text{ is unimodal, symmetric about the} \\ &\quad \text{mode (or perhaps median) of } w_0\} \\ \mathcal{V} &= \{v = (1 - \epsilon_2)v_0 + \epsilon_2q_v: q_v \text{ is unimodal, symmetric about the} \\ &\quad \text{mode (or median) of } v_0\}, \end{aligned}$$

and apply the ideas of Section 3. We do not attempt the analysis here, because nothing new conceptually is involved and the argument would be moderately lengthy.

5.3. *First and second stage uncertainty.* The simplest modification of (5.7) that introduces uncertainty in the first stage of the prior is simply to add an arbitrary overall contamination. Thus we consider

$$\pi(\theta) = (1 - \epsilon_1 - \epsilon_2)\pi_0 + \epsilon_1q_1 + \epsilon_2q_2,$$

where  $q_2 \in \mathcal{L}_2 = \mathcal{P}$ ,

$$q_1 = \int \left[ \prod_{i=1}^p g_{\omega}(\theta_i) \right] s(\omega) d\omega \in \mathcal{L}_1 = \{q: s \in \mathcal{S}\},$$

and  $\pi_0$ ,  $s$ , and  $\mathcal{S}$  are as in (5.7). In other words,  $q_1$  arises from possible second stage prior uncertainty, while  $q_2$  allows for basic error in the empirical Bayes model.

Allowing arbitrary  $q_2$  is again, probably excessively crude. In particular, complete abandonment of the empirical Bayes structure may be unrealistic. For illustrative purposes, however, this is convenient.

As mentioned in Section 5.1, the ML-II prior can be found (here, at least) by separately maximizing  $m(x|q_1)$  and  $m(x|q_2)$ . Maximization of  $m(x|q_1)$  was discussed in the previous section. And  $m(x|q_2)$  will simply be maximized when  $q_2$  is a unit point mass at  $\hat{\theta}$ , the maximum likelihood estimate. Thus the ML-II prior is [assuming  $\mathcal{S} = \mathcal{P}$  and letting  $I(\hat{\theta})$  denote a unit point mass at  $\hat{\theta}$ ]

$$\hat{\pi}(\theta) = (1 - \epsilon_1 - \epsilon_2)\pi_0(\theta) + \epsilon_1 \left[ \prod_{i=1}^p g_{\omega}(\theta_i) \right] + \epsilon_2 I(\hat{\theta}).$$

Formulas (5.3)–(5.5) can now easily be employed to give desired conclusions. In the situation of Example 3, for instance, all calculations can be carried out explicitly; indeed, the needed modifications to the formulae there are very minor and so will be omitted. The behavior of  $\delta^{\hat{\pi}}$ , the ML-II posterior mean, is worth

mentioning, however. If the data are compatible with  $\pi_0$  (i.e., are near  $\mu_0$ ) then the conjugate prior posterior mean  $\delta^0$  will dominate; if the data are similar but not near  $\mu_0$ , then  $\delta^{\hat{\pi}}$  will be close to the natural empirical Bayes rule  $\hat{\delta}^a$ ; and if the data are not compatible with the empirical Bayes model, then  $\delta^{\hat{\pi}}$  will be close to the maximum likelihood estimate,  $\hat{\theta} = x$ .

**6. Discussion.** We view this paper as a hopeful first step in the development of systematic robust Bayesian analyses for rich classes of priors. The original goals of the paper were (i) to demonstrate that it is possible to work with complex classes of priors (as in Section 4), and to indicate mathematical techniques for doing so; (ii) to point out the numerous intuitive and calculational reasons for approaching robustness through consideration of  $\varepsilon$ -contamination classes; and (iii) to exhibit the value of the ML-II approach in obtaining "robust priors." Through tenacious prodding from skeptical referees and the associate editor, we have become more cautious in our assessment of success in goal (iii). A brief discussion of this issue is in order.

The key to obtaining a robust prior appears to be the selection of a prior with tails that are much flatter than the tails of the likelihood function [see Berger (1984, 1985) for discussion and references]. Unfortunately, this observation does not provide a readily implementable "solution" to robustness questions. Basic difficulties include (i) the uncertainty as to the choice of robust prior tails and (ii) the calculational complexities that can result. In addition to the already established computational simplicity, our hope for the ML-II technique was that it would *automatically* provide a "prior" robust against the type of deviations considered plausible.

It is important to explain our reasons for believing that ML-II would succeed when applied to  $\varepsilon$ -contamination classes. Suppose that  $\varepsilon$  is fairly small and that  $\Gamma$  contains all plausible priors, but none that are terribly implausible. Consider first the case where  $m(x|\pi_0)$  is large, i.e., the data is compatible with the nominally specified  $\pi_0$ . This is a situation of nearly automatic robustness, in that, since the central portions of all  $\pi$  in  $\Gamma$  (and, hence,  $\hat{\pi}$ ) will be similar to that of  $\pi_0$ , the conclusions will be very similar for all  $\pi$  in  $\Gamma$ . On the other hand, consider the case where  $x$  is such that  $m(x|\pi_0)$  is small. This is precisely the situation in which the prior tail is highly influential and the use of a large prior tail is desirable. The ML-II technique will naturally select a prior with a large tail, since such priors are those most compatible with the data. Opposing the above encouraging tendencies toward robustness, is the danger that data-selected priors will tend to over-concentrate about the likelihood function, thereby yielding error estimates that are too small. (Such dangers lurk in the shadows of much of empirical Bayesian analysis and for that matter the use of data-selected models.) When  $\Gamma$  contains unreasonable priors, as in Section 2, this problem can completely dominate. However, the hope was that reasonable  $\Gamma$ , such as those in Sections 3 and 4, by limiting the possible concentration about the likelihood function, would not succumb to this danger to a serious extent.

To examine the degree to which this hope was realized, we return to the example discussed in Sections 2-4:  $X \sim \mathcal{N}(\theta, 1)$ ,  $\pi_0$  is  $\mathcal{N}(0, 2)$ , and  $\varepsilon = 0.2$ . The

four estimators of  $\theta$  presented were (changing notation for convenience):

- $\delta_1$ , the Bayes estimator with respect to  $\pi_0$ ;
- $\delta_2$ , the ML-II posterior mean for  $\mathcal{Q}_2 = \{\text{all distributions}\}$ ;
- $\delta_3$ , the ML-II posterior mean for  $\mathcal{Q}_3 = \{\text{all symmetric, unimodal distributions}\}$ ;
- $\delta_4$ , the ML-II posterior mean for  $\mathcal{Q}_4 = \{q \text{ such that } \pi \text{ is unimodal}\}$ .

Let us add to this list  $\delta_5$ , the posterior mean for a  $t$ -prior distribution with 4 degrees of freedom, median zero, and quartiles equal to  $\pm 0.96$  (the quartiles of  $\pi_0$ ); this corresponds to a scale factor of 1.3 for the  $t$ -distribution. The point of including the  $t$ -prior is that it is a reasonable robust prior (which happens to be in all the classes  $\Gamma$  for  $\mathcal{Q}_2$ ,  $\mathcal{Q}_3$ , and  $\mathcal{Q}_4$ ), and provides a benchmark for judging the

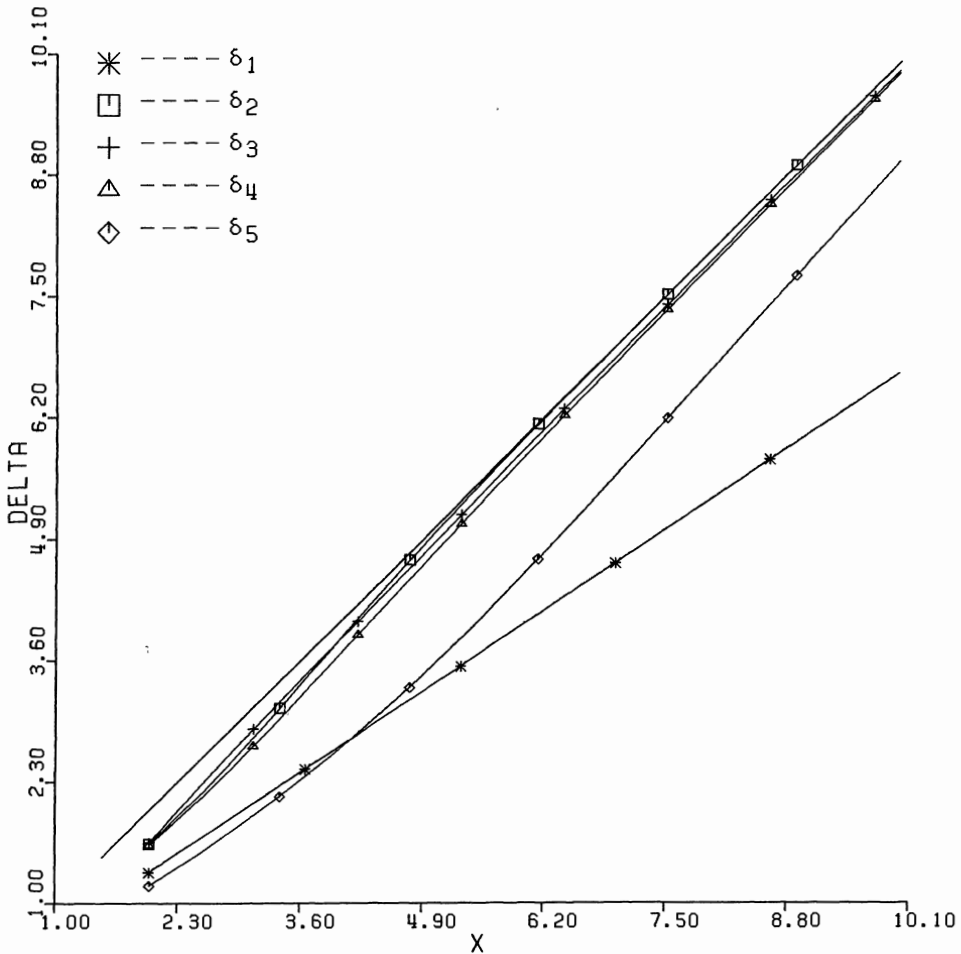


FIG. 1.



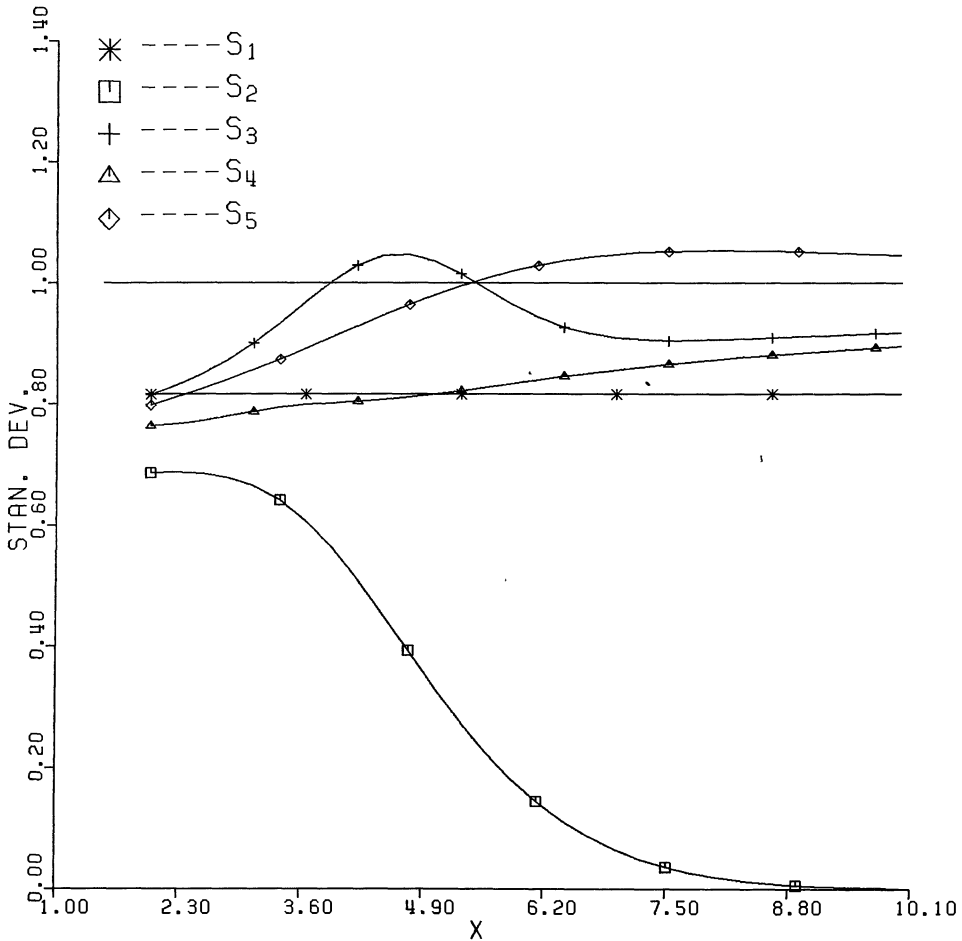


FIG. 2.

performance of ML-II. Figure 1 presents graphs of  $\delta_1$  through  $\delta_5$  for  $2 \leq x \leq 10$  (there is very little difference between the estimators for  $x < 2$ ). Figure 2 presents the corresponding natural error estimates, namely, the posterior (ML-II versions for  $\delta_2, \delta_3,$  and  $\delta_4$ ) standard deviations  $s_1(x), s_2(x), \dots, s_5(x)$ .

The most obvious conclusion from Figure 1 is that the ML-II estimates are more conservative (closer to  $x$ ) than the Bayes estimates. Also, quite naturally, larger  $\mathcal{L}$  result in more conservative estimates (note that  $\mathcal{L}_3 \subset \mathcal{L}_4 \subset \mathcal{L}_2$ ).

Turning to Figure 2, note first that  $s_2(x)$  is indeed fairly ridiculous for large  $x$ . Also,  $s_4(x)$  is moderately smaller than  $s_5(x)$ . We had hoped that the attractive  $\Gamma$  of Section 4 would not be so large as to result in an excessively small ML-II error estimate, but when compared to  $s_5$ , the behavior of  $s_4$  is borderline. On the other hand, the behavior of  $s_3$  is definitely satisfactory. Finally, note that for moderate  $x$ , both  $s_3$  and  $s_5$  rise above the conditional standard deviation (1) of  $X$ . [The

behavior of  $s_5$  is in fact quite general for robust priors; for example, see O'Hagan, (1981)].

To summarize, we proposed the ML-II technique as a possible automatic "robustifier" of standard (conjugate) priors. We have shown that the ML-II technique is computationally feasible, and that it can successfully robustify  $\pi_0$  if the class  $\Gamma$  is sensible (i.e., does not contain silly contaminations such as those in  $\mathcal{Q}_2$ ). Alternatively, in any given problem one could (should?) attempt to construct a robust prior, such as the  $t$ -prior in the above example. We are in no way opposed to such efforts, but there are numerous technical and theoretical issues involved in insuring that robustness is obtained; the process is far from automatic. While we are not strong proponents of "automation" in statistics (few Bayesians are), we recognize the forces driving statistics in that direction. Of course, nothing is completely automatic; our ML-II approach does require the imputation of  $\epsilon$  and  $\mathcal{Q}$ . However, it will often be reasonable to choose  $\mathcal{Q}$  in a standard or default fashion, say as in Section 3. The quantity  $\epsilon$  is reasonably accessible to intuition, relating in a fairly straightforward manner to one's confidence in the specification of  $\pi_0$ .

**Acknowledgments.** We would like to thank the referees and Associate Editor for very valuable comments and suggestions.

## REFERENCES

- BERGER, J. (1982). Bayesian robustness and the Stein effect. *J. Amer. Statist. Assoc.* **77** 358–368.
- BERGER, J. (1984). The robust Bayesian viewpoint (with Discussion). In *Robustness in Bayesian Statistics* (J. Kadane, ed.). North-Holland, Amsterdam.
- BERGER, J. (1985). *Statistical Decision Theory and Bayesian Analysis*. Springer, New York.
- BERGER, J. and BERLINER, M. (1983). Robust Bayes and empirical Bayes analysis with  $\epsilon$ -contaminated priors. Technical Report 83-35, Purdue Univ., West Lafayette, Indiana.
- BERGER, J. and BERLINER, M. (1984). Bayesian input in Stein estimation and a new minimax empirical Bayes estimator. *J. Econometrics* **25** 87–108.
- BERGER, J. and SELLEKE, T. (1984). Testing a point null hypothesis: the irreconcilability of significance levels and evidence. Technical Report 84-27, Purdue Univ., West Lafayette, Indiana.
- BICKEL, P. J. (1984). Parametric robustness or small biases can be worthwhile. *Ann. Statist.* **12** 864–879.
- BISHOP, Y. M. M., FIENBERG, S. E. and HOLLAND, P. W. (1975). *Discrete Multivariate Analysis*. M.I.T. Press, Cambridge.
- BLUM, J. R. and ROSENBLATT, J. (1967). On partial a priori information in statistical inference. *Ann. Math. Statist.* **38** 1671–1678.
- BOX, G. E. P. (1980). Sampling and Bayes' inference in scientific modelling and robustness (with Discussion). *J. Roy. Statist. Soc. Ser. A* **143** 383–430.
- BOX, G. E. P. and TIAO, G. C. (1973). *Bayesian Inference in Statistical Analysis*. Addison-Wesley, Reading, Mass.
- DEMPSTER, A. P. (1975). A subjectivist look at robustness. *Bull. Inst. Internat. Statist.* **46** 349–374.
- DEROBERTIS, L. and HARTIGAN, J. A. (1981). Bayesian inference using intervals of measures. *Ann. Statist.* **9** 235–244.
- GOOD, I. J. (1950). *Probability and the Weighing of Evidence*. Griffin, London.
- GOOD, I. J. (1965). *The Estimation of Probabilities*. M.I.T. Press, Cambridge.
- GOOD, I. J. (1980). Some history of the hierarchical Bayesian methodology. In *Bayesian Statistics* (J. M. Bernardo, M. H. DeGroot, D. V. Lindley, and A. F. M. Smith, eds.). University Press, Valencia.

- HUBER, P. J. (1973). The use of Choquet capacities in statistics. *Bull. Inst. Internat. Statist.* **45** 181–191.
- JEFFREYS, H. (1961). *Theory of Probability*. 3rd ed. University Press, Oxford.
- KADANE, J. B., DICKEY, J. M., WINKLER, R. L., SMITH, W. S. and PETERS, S. C. (1980). Interactive elicitation of opinion for a normal linear model. *J. Amer. Statist. Assoc.* **75** 845–854.
- LEAMER, E. E. (1978). *Specification Searches*. Wiley, New York.
- LEONARD, T. (1974). A modification to the Bayes estimate for the mean of a normal distribution. *Biometrika* **61** 627–628.
- LINDLEY, D. V. (1971). The estimation of many parameters. In *Foundations of Statistical Inference* (V. P. Godambe and D. A. Sprott, eds.). Holt, Rinehart, and Winston, Toronto.
- LINDLEY, D. V. and SMITH, A. F. M. (1972). Bayes estimates for the linear model. *J. Roy. Statist. Soc. Ser. B* **34** 1–41.
- MARAZZI, A. (1985). Robust Bayesian estimation for the linear model. *Statistics and Decisions* **3**.
- MARITZ, J. S. (1970). *Empirical Bayes Methods*. Methuen, London.
- MORRIS, C. (1983). Parametric empirical Bayes inference: theory and applications (with Discussion). *J. Amer. Statist. Assoc.* **78** 47–65.
- O'HAGAN, A. (1981). A moment of indecision. *Biometrika* **68** 329–330.
- SCHNEEWEISS, H. (1964). Eine Entscheidungsregel für den Fall partiell bekannter Wahrscheinlichkeiten. *Unternehmensforschung* **8** 86–95.
- ZELLNER, A. (1985). On assessing prior distributions and Bayesian regression analysis with  $g$ -prior distributions. In *Bayesian Inference and Decision Techniques with Applications* (P. K. Goel and A. Zellner, eds.). North-Holland, Amsterdam.

DEPARTMENT OF STATISTICS  
PURDUE UNIVERSITY  
WEST LAFAYETTE, INDIANA 47907

DEPARTMENT OF STATISTICS  
OHIO STATE UNIVERSITY  
1958 NEIL AVENUE  
COLUMBUS, OHIO 43210