

THE STEIN EFFECT AND BAYESIAN ANALYSIS:  
A REEXAMINATION

by

Jean-Francois Angers and James O. Berger  
Purdue University

Technical Report #85-6

Department of Statistics  
Purdue University

May 1985

THE STEIN EFFECT AND BAYESIAN ANALYSIS:  
A REEXAMINATION

Jean-Francois Angers and James O. Berger

Statistics Department, Purdue University, West Lafayette, IN 47907, U.S.A.

ABSTRACT

The Stein effect, that one could improve frequentist risk by combining "independent" problems, has long been an intriguing paradox to statistics. We briefly review the Bayesian view of the paradox, and indicate that previous justifications of the Stein effect, through concerns of "Bayesian robustness," were misleading. In the course of doing so, several existing robust Bayesian and Stein-effect estimators are compared for a variety of situations.

Key Words and Phrases: Stein effect; Bayesian robustness; independent problems; classes of priors.

## 1. INTRODUCTION

Suppose  $X_i \sim \mathcal{N}(\theta_i, \sigma_i^2)$  ( $\sigma_i^2$  known), and that it is desired to estimate  $\theta_i$  by an estimator  $\delta_i(x_i)$  under squared error loss. The standard estimator  $\delta_i^0(x_i) = x_i$  is admissible when frequentist risk (expected loss) is the criterion. Stein (1955) showed, however, that if one combined  $p$  such problems, i.e., considered the problem of estimating  $\theta = (\theta_1, \dots, \theta_p)^t$ , based on  $X = (X_1, \dots, X_p)^t$  under sum of squares error loss, then the corresponding natural estimator,  $\delta^0(X) = (\delta_1^0(x_1), \dots, \delta_p^0(x_p))^t = X$ , is inadmissible when  $p \geq 3$ . Indeed, James and Stein (1960) showed that (when all  $\sigma_i^2 = 1$ ) a better estimator is given by

$$\delta^{J-S}(X) = \left(1 - \frac{(p-2)}{|X|^2}\right) X; \quad (1.1)$$

thus, in terms of frequentist risk

$$R(\theta, \delta) = E_{\theta} |\delta(X) - \theta|^2,$$

James and Stein (1960) showed that  $R(\theta, \delta^{J-S}) < R(\theta, \delta^0)$  for all  $\theta$ .

Note that  $\delta^{J-S}$  involves all the  $x_j$  in estimating each  $\theta_i$ . This is counter-intuitive when the  $\theta_i$  are from completely different problems, as in the following example.

Example 1. Suppose  $\theta_1$  is the mean corn yield per acre in Indiana in 1985,  $\theta_2$  is the mass of the universe, and  $\theta_3$  is the age of a newly discovered Greek artifact. A large business conglomerate is simultaneously interested in estimating these three quantities, and the directors ascertain that the overall loss to the business in mis-estimation is  $[(\theta_1 - \delta_1)^2 + (\theta_2 - \delta_2)^2 + (\theta_3 - \delta_3)^2]$ . It

strikes most people as rather odd to combine measurements involving such  $\theta_j$  in a manner such as (1.1).

This phenomenon, that one can combine independent problems and obtain improved estimators, is called the Stein effect. It is a very general phenomenon in frequentist decision theory (see Berger (1985a) for references). Several explanations for the phenomenon have been proposed. One common explanation (cf. Efron and Morris (1973)) is that if the  $\theta_j$  are related, as in being i.i.d. observations from a common population, then estimators similar to  $\delta_{\hat{\theta}}^{J-S}$  can arise as empirical Bayes estimators. However, while such situations are extremely important areas for application of Stein estimation, they are not cases in which the coordinate problems are truly independent; since the  $\theta_j$  arise from a common population, knowledge concerning  $\theta_j$  is clearly relevant to estimating  $\theta_i$ . We will be considering instead situations such as Example 1, where the Stein effect holds, but in which it would be absurd to think of the  $\theta_j$  as coming from a common population.

A second common "explanation" of the Stein effect is that it is natural because the losses for the different problems have been combined. Indeed, such combination can force a dependence into the final estimate when one has a non-linear utility function. The Stein phenomenon is usually stated, however, in terms of combining the individual losses,  $L_i(\theta_i, \delta_i)$ , in a linear fashion, to yield

$$L(\hat{\theta}, \hat{\delta}) = \sum_{i=1}^p L_i(\theta_i, \delta_i) \quad (1.2)$$

(or perhaps some weighted linear combination), and the situation is then not so clear. It may seem that this forcing of the  $L_i$  to be on a common scale might be the cause of a dependence in estimation, but a number of intuitive arguments refute this. The simplest, which we will pursue in this paper, is the Bayesian view of the situation.

Let  $\pi(\theta)$  denote a prior (density for convenience) for  $\theta$ . The assumption that the problems are completely unrelated can be expressed by the Bayesian condition that the  $\theta_i$  are a priori independent, i.e., that

$$\pi(\theta) = \prod_{i=1}^p \pi_i(\theta_i). \quad (1.3)$$

Now, for such a prior and with a loss as in (1.2), it is straightforward to verify that the Bayes rule is

$$\delta^{\pi}(x) = (\delta^{\pi_1}(x_1), \dots, \delta^{\pi_p}(x_p))^t,$$

where  $\delta^{\pi_i}(x_i)$  is the Bayes rule for estimating  $\theta_i$  under loss  $L_i$  and with respect to the prior  $\pi_i$ . Thus a Bayesian feels that, if the problems are truly unrelated (and the prior is known), then only  $x_i$  should be involved in the estimation of  $\theta_i$ . The combination of losses in a linear fashion, such as (1.2), has no effect on the Bayes rule.

To obtain a Bayesian explanation for the Stein effect in a setting such as Example 1, we must, therefore, leave the pure Bayesian framework. A natural and practically relevant direction to look is towards concerns of Bayesian robustness. Indeed, Berger (1982a) so approached the problem, imagining the specification

of a prior  $\pi^0$ , of the form (1.3), and advocating consideration of the class of priors

$$\Gamma = \{\pi = (1-\epsilon) \pi^0 + \epsilon q; q \text{ an arbitrary distribution}\}, \quad (1.4)$$

where  $0 < \epsilon < 1$  is a constant reflecting the accuracy of the specification of  $\pi^0$ . The idea is that, in practice,  $\pi^0$  will merely be an approximate quantification of prior beliefs, and that priors "close to"  $\pi^0$  are also plausible. The robust Bayesian seeks to conduct an analysis which is likely to be satisfactory for all such "close" priors. The class in (1.4) is a simple and useful class of close priors. (It might be argued that allowing the "contamination,"  $q$ , to be arbitrary allows too much deviation from  $\pi^0$ , but for the purposes of this paper such a concern is irrelevant; the same results can be shown to hold as long as there are priors in  $\Gamma$  corresponding to  $q$  which give arbitrarily small mass to compact sets).

In Berger (1982a, 1984), it was shown that the Stein effect is valuable in achieving Bayesian robustness with respect to classes such as (1.4); indeed, it was shown to be possible (in symmetric situations) to construct Stein-type estimators which are nearly optimal in terms of Bayes risk

$$r(\pi, \delta) = E^\pi R(\theta, \delta) = \int R(\theta, \delta) \pi(\theta) d\theta, \quad (1.5)$$

for an elicited prior  $\pi^0$ , and yet are also highly robust with respect to  $\Gamma$ , in the sense that

$$r_\Gamma(\delta) = \sup_{\pi \in \Gamma} r(\pi, \delta) \quad (1.6)$$

is also nearly optimal. The criteria (1.5) and (1.6) are, of course, partly frequentist in nature, and are of interest to a Bayesian only in that they can indicate good average posterior expected loss for  $\delta$ . (In Berger (1984, 1985b) the value of such criteria to a Bayesian is more fully discussed.) It suffices here to note that these criteria are of interest, and that Berger (1982a, 1982b) showed that a Stein-type estimator can do much better by these criteria than can robust coordinatewise independent estimators of the form

$$\delta(x) = (\delta_1(x_1), \dots, \delta_p(x_p))^t. \quad (1.7)$$

Thus there appeared to be value in the Stein effect, even from a Bayesian perspective.

The intuitive paradox has caused us to continue to examine the issue, and we noticed a disturbing step in the above argument. The step was the inclusion of priors in (1.4) that are coordinatewise dependent; even though  $\pi^0$  is of the form (1.3), many (indeed most) of the priors in (1.4) will have dependent coordinates. And, although it may often be the case that dependence among the  $\theta_i$  is possible (in which case (1.4) might be reasonable), our stated goal was to investigate the value of the Stein effect for completely unrelated problems. Thus, in Example 1, while we might be quite uncertain as to the choice of  $\pi^0$ , our prior belief in the independence of the  $\pi_i$  may be so strong that only priors satisfying (1.3) would be deemed reasonable.

A sensible class of priors for such a situation is

$$\Gamma_1 = \left\{ \pi = \prod_{i=1}^p [(1-\varepsilon) \pi_i^0(\theta_i) + \varepsilon q_i(\theta_i)], \text{ the } q_i \text{ arbitrary} \right\}. \quad (1.8)$$

(Again, the only essential requirement for the following results is that  $\Gamma_1$  contain priors corresponding to  $q_i$  which give arbitrarily small mass to compact sets.) This class contains priors which are close to  $\pi^0$  and which preserve the independence structure. The obvious question to ask is--Does the Stein effect allow improved performance, from a robust Bayesian perspective, for  $\Gamma_1$ ? We will see that the answer is, essentially, no; there does not appear to be value to the Stein effect in combining independent problems.

To avoid possible confusion, we again emphasize that we are not considering the formally related empirical Bayes situation where  $\pi$  is of the form (1.3), and the coordinatewise priors,  $\pi_i$ , have unknown and partially common features. The coordinate problems are not then independent to a Bayesian (placing a second stage prior on the common features and integrating out clearly reveals the dependence), and a Bayesian will naturally employ coordinatewise dependent estimators.

## 2. PRELIMINARIES

For convenience, we restrict ourselves to the symmetric situation with independent  $X_i \sim \mathcal{N}(\theta_i, \sigma^2)$ ,  $i=1, \dots, p$ ,  $\sigma^2$  known, loss  $L_i(\theta_i, \delta_i) = (\theta_i - \delta_i)^2$  in estimating  $\theta_i$  by  $\delta_i$ , and overall loss  $L(\theta, \delta) = \sum_{i=1}^p (\theta_i - \delta_i)^2$  in estimating  $\theta = (\theta_1, \dots, \theta_p)^t$  by  $\delta = (\delta_1, \dots, \delta_p)^t$ . We imagine that the  $\theta_i$  are a priori independent, and that rough quantification of the prior information results in a  $\mathcal{N}(\mu_i, \tau^2)$  prior, to be denoted  $\pi_i^0$ , for each  $\theta_i$ . Here we presume that  $\tau^2$  and  $\mu = (\mu_1, \dots, \mu_p)^t$  are subjectively specified numbers. Thus the "base" prior for the Bayesian robustness investigation is  $\pi^0(\theta) = \prod_{i=1}^p \pi_i^0(\theta_i)$ , the  $\mathcal{N}_p(\mu, \tau^2 I)$  distribution. Note that, since  $\tau^2$  is assumed known, this is not an empirical Bayes situation. It is a situation where, by coincidence, we have  $p$  completely independent problems which have similar known variances. This is not being put forth as a realistic practical situation (although the amount of theoretical



literature on the situation is quite astounding), but rather because it is the situation most conducive to improvements via the Stein effect; if the Stein effect is of no or limited value to Bayesian robustness here, it is unlikely to be helpful in more realistic scenarios involving completely independent problems with very different variances.

We will work with two classes of prior distributions. The first has already been discussed, namely  $\Gamma_1$  in (1.8). The second class is somewhat less conservative, allowing at most one of the coordinates of  $\theta$  to have a contaminated prior. The class is

$$\Gamma_2 = \{ \pi(\theta) = \frac{1}{p} \sum_{j=1}^p ([(1-\epsilon)\pi_j^0(\theta_j) + \epsilon q_j(\theta_j)] \prod_{i \neq j} \pi_i^0(\theta_i)) \}, \text{ the } q_j \text{ arbitrary}. \quad (2.1)$$

Note that  $\Gamma_2 \subset \Gamma_1$ , so that demanding robustness with respect to  $\Gamma_2$  is less stringent than demanding it for  $\Gamma_1$ .

By looking at  $r(\pi^0, \delta)$ ,  $r_{\Gamma_1}(\delta)$ , and  $r_{\Gamma_2}(\delta)$  (see (1.5) and (1.6)), we hope to obtain reasonable indications of the performance of an estimator  $\delta$ . It is actually somewhat more convenient to normalize  $r(\pi^0, \delta)$  and  $r_{\Gamma}(\delta)$  as follows: define

$$\text{RSR}(\pi^0, \delta) = \frac{r(\pi^0, \delta) - r(\pi^0, \delta^{\pi^0})}{r(\pi^0, \delta^0) - r(\pi^0, \delta^{\pi^0})} = \frac{(\sigma^2 + \tau^2)}{p\sigma^4} r(\pi^0, \delta) - \frac{\tau^2}{\sigma^2}, \quad (2.2)$$

and

$$r_{\Gamma}^*(\delta) = \frac{(\sigma^2 + \tau^2)}{p\sigma^4} r_{\Gamma}(\delta) - \frac{\tau^2}{\sigma^2}. \quad (2.3)$$

The "relative savings risk" of  $\delta$ ,  $RSR(\pi^0, \delta)$ , measures the proportion of available improvement over  $\delta^0(x) = x$  that is achieved by  $\delta$  (for the prior  $\pi^0$ ).  $RSR$  near 0 is optimal, while  $RSR$  near 1 indicates performance as bad as that of  $\delta^0$ . (Efron and Morris (1971) first introduced this concept.) The main purpose in using these scaled versions of risk is that the scaled versions are often independent of  $\sigma^2$ ,  $\tau^2$ , and  $\mu$ .

Note also that  $r_\Gamma$  (or  $r_\Gamma^*$ ) can be used to define  $\Gamma$ -minimaxity. A  $\Gamma$ -minimax estimator is an estimator which minimizes  $r_\Gamma$  (or, equivalently,  $r_\Gamma^*$ ) over all possible estimators.

### 3. THE ESTIMATORS AND RISKS

The estimators that are to be considered are various robustified versions of the Bayes estimator with respect to  $\pi^0$ , which is given for the situation of Section 2 by

$$\delta^{\pi^0}(x) = (\delta_1^{\pi^0}(x_1), \dots, \delta_p^{\pi^0}(x_p))^T,$$

where

$$\delta_i^{\pi^0}(x_i) = x_i - \frac{\sigma^2}{(\sigma^2 + \tau^2)} (x_i - \mu_i).$$

Clearly  $RSR(\pi^0, \delta^{\pi^0}) = 0$ , but it is easy to see that  $r_{\Gamma_1}^*(\delta^{\pi^0}) = r_{\Gamma_2}^*(\delta^{\pi^0}) = \infty$ , so that  $\delta^{\pi^0}$  is highly suspect from a robustness viewpoint. (Note that  $r_\Gamma^*(\delta^0) = 1$  for any  $\Gamma$ , so that values of  $r_\Gamma^*$  larger than one can be avoided.)

The most natural way to robustify  $\hat{\delta}_\pi^0$  is on a coordinatewise basis. An appealing choice is the limited translation estimator of Efron and Morris (1971), given by

$$\hat{\delta}_\pi^{M,I}(\underline{x}) = (\delta_1^{M,I}(x_1), \dots, \delta_p^{M,I}(x_p))^t, \quad (2.4)$$

where

$$\delta_i^{M,I}(x_i) = x_i - \min \left\{ 1, \frac{[M(\sigma^2 + \tau^2)]^{\frac{1}{2}}}{|x_i - \mu_i|} \right\} \frac{\sigma}{(\sigma^2 + \tau^2)} (x_i - \mu_i).$$

A robustified version of  $\hat{\delta}_\pi^0$  that involves the Stein effect is given in Berger (1982b), and can be written as

$$\hat{\delta}_\pi^{M,S}(\underline{x}) = \underline{x} - \min \{1, \rho_{M,p}(|\underline{x} - \underline{\mu}|^2 / (\sigma^2 + \tau^2))\} \frac{\sigma^2}{(\sigma^2 + \tau^2)} (\underline{x} - \underline{\mu}), \quad (2.5)$$

where, for  $p = 3, 5$ , and  $7$  (the 3 cases to be considered in this paper),  $\rho_{M,p}(r)$  is given by

$$\rho_{M,3}(r) = \frac{\sqrt{M}}{\sqrt{r}} + \frac{2}{r}, \quad \rho_{M,5}(r) = \frac{M}{2 + \sqrt{Mr}} + \frac{6}{r},$$

$$\rho_{M,7}(r) = \frac{M(2 + \sqrt{Mr})}{12 + 6\sqrt{Mr} + Mr} + \frac{10}{r}.$$

For general discussion of the properties of these estimators, see Berger (1982b). Their Bayesian robustness is indicated by the fact that they equal  $\delta_\pi^0$  if  $|\bar{x} - \mu|$  is small (i.e., the prior information appears to be correct), and yet are not allowed to deviate excessively from  $\delta_\pi^0(\bar{x}) = \bar{x}$  if  $|\bar{x} - \mu|$  is large.

For various  $\varepsilon$  and  $\Gamma_1$  or  $\Gamma_2$  (defined in (1.8) or (2.1), respectively), we will calculate and compare the "most robust" estimators of the forms (2.4) and (2.5). This will be done by simply minimizing  $r_{\Gamma}^*(\delta_\pi^{M,I})$  or  $r_{\Gamma}^*(\delta_\pi^{M,S})$  over  $M$  (for the various  $\varepsilon$  and  $\Gamma$ ). For notational convenience,  $M^*$  will be used to denote the optimal value of  $M$  in each situation. (Note that this will, in general, depend on the choice of  $\Gamma$ ,  $\varepsilon$ , and the form of the estimator.)

The interest in  $\delta_\pi^{M^*,I}$  is that it can be shown to be nearly optimal among all coordinatewise independent estimators (i.e., those of the form (1.7)), in terms of minimizing  $r_{\Gamma_1}^*$  or  $r_{\Gamma_2}^*$ . Likewise,  $\delta_\pi^{M^*,S}$  is nearly optimal among all "Stein-effect" estimators of the form  $\delta_\pi(\bar{x}) = \bar{x} - \rho(|\bar{x} - \mu|^2)(\bar{x} - \mu)$ . (The arguments are similar to those in Berger and Berliner (1984) or Marazzi (1985), and will be omitted.) Thus a comparison of the performance of these estimators should provide a strong indication as to the value of the Stein effect in the situations considered. The following lemmas present risk formulas needed for the calculations and comparisons.

Lemma 1. The estimator  $\delta_\pi^{M,I}$  has

$$RSR(\pi^0, \delta_\pi^{M,I}) = [1 - \psi_1(M)][1 + M] - \left(\frac{2M}{\pi}\right)^{\frac{1}{2}} e^{-M/2}, \quad (2.6)$$

where  $\psi_1$  is the c.d.f. of the chi-square distribution with one degree of freedom.

Also

$$r_{\Gamma_1}^*(\hat{\delta}^{M,I}) = (1 - \varepsilon)RSR(\pi^0, \hat{\delta}^{M,I}) + \varepsilon(M + 1) \quad (2.7)$$

and

$$r_{\Gamma_2}^*(\hat{\delta}^{M,I}) = (1 - \frac{\varepsilon}{p})RSR(\pi^0, \hat{\delta}^{M,I}) + \frac{\varepsilon}{p}(M + 1). \quad (2.8)$$

Proof. Equation (2.6) is given in a slightly different form in Efron and Norris (1971). To verify (2.7), observe that, since  $\Gamma_1$  contains only priors of the form (1.3),

$$\begin{aligned} \sup_{\pi \in \Gamma_1} r(\pi^0, \hat{\delta}^{M,I}) &= \sup_{\pi \in \Gamma_1} \sum_{i=1}^p r(\pi_i, \delta_i^{M,I}) \\ &= \sup_{\pi \in \Gamma_1} \sum_{i=1}^p (1 - \varepsilon)r(\pi_i^0, \delta_i^{M,I}) + \varepsilon r(q_i, \delta_i^{M,I}) \\ &= (1 - \varepsilon)r(\pi^0, \hat{\delta}^{M,I}) + \varepsilon \sum_{i=1}^p \sup_{q_i} r(q_i, \delta_i^{M,I}) \\ &= (1 - \varepsilon)r(\pi^0, \hat{\delta}^{M,I}) + \varepsilon \sum_{i=1}^p \sup_{\theta_i} R(\theta_i, \delta_i^{M,I}) \\ &= (1 - \varepsilon)r(\pi^0, \hat{\delta}^{M,I}) + \varepsilon p \left( \frac{M\sigma^4}{(\sigma^2 + \tau^2)} + \sigma^2 \right), \end{aligned}$$

the last fact following from a result of Efron and Morris (1971). Using the

definitions of RSR and  $r_{\Gamma_1}^*$  and simplifying, yields (2.7). Equation (2.8) is similarly established.  $\square$

Lemma 2. The estimator  $\delta_{\zeta}^{M,S}$  has

$$\text{RSR}(\pi^0, \delta_{\zeta}^{M,S}) = [1 - \psi_p(y)][1 + \frac{M}{p}] - \frac{(y/2)^{p/2} e^{-y/2}}{\Gamma(1 + p/2)}, \quad (2.9)$$

where  $\psi_p$  is the c.d.f. of the chi-square distribution with  $p$  degrees of freedom and  $y$  is the solution to the equation

$$\rho_{M,p}(y) = 1.$$

Also,

$$r_{\Gamma_1}^*(\delta_{\zeta}^{M,S}) = (1 - \epsilon)^p \text{RSR}(\pi^0, \delta_{\zeta}^{M,S}) + [1 - (1 - \epsilon)^p](\frac{M}{p} + 1), \quad (2.10)$$

and

$$r_{\Gamma_2}^*(\delta_{\zeta}^{M,S}) = (1 - \epsilon) \text{RSR}(\pi^0, \delta_{\zeta}^{M,S}) + \epsilon(\frac{M}{p} + 1). \quad (2.11)$$

Proof. Equation (2.9) follows from Berger (1982b) or Chen (1983). To establish (2.10), note first that any distribution of the form

$$\pi^*(\theta) = [\prod_{\substack{\text{some} \\ \text{indices}}} \pi_i^0(\theta_i)] [\prod_{\substack{\text{other} \\ \text{indices}}} q_i(\theta_i)] \quad (2.12)$$

satisfies

$$\begin{aligned} \sup_{\{q_i\}} r(\pi^*, \delta^M, S) &= \sup_{\theta} R(\theta, \delta^M) \\ &= p\sigma^2 + \frac{M\sigma^4}{(\sigma^2 + \tau^2)}. \end{aligned}$$

(These equalities follow by letting one of the involved  $q_i$  give mass to a point  $\theta_i^*$ , letting  $|\theta_i^*| \rightarrow \infty$ , and using results about the behavior of  $\rho_{M,p}$  from Berger (1982b).) But any  $\pi \in \Gamma_1$  can be written

$$\pi(\theta) = (1 - \varepsilon)^p \pi^0(\theta) + \sum_{\ell} \pi_{\ell}^*(\theta),$$

where the  $\pi_{\ell}^*$  are of the form (2.12) and have total mass  $[1 - (1 - \varepsilon)^p]$ ; thus

$$\begin{aligned} \sup_{\pi \in \Gamma_1} r(\pi, \delta^M, S) &= (1 - \varepsilon)^{\bar{p}} r(\pi^0, \delta^M, S) + \sup_{\{q_i\}} \sum_{\ell} r(\pi_{\ell}^*, \delta^M, S) \\ &= (1 - \varepsilon)^p r(\pi^0, \delta^M, S) + [1 - (1 - \varepsilon)^p] \left[ p\sigma^2 + \frac{M\sigma^4}{(\sigma^2 + \tau^2)} \right]. \end{aligned}$$

Using the definitions of RSR and  $r_{\Gamma_1}^*$  and simplifying, yields (2.10). The verification of (2.11) is similar and will be omitted.  $\square$

Having derived explicit formulas for  $r_{\Gamma_1}^*$  and  $r_{\Gamma_2}^*$  for the estimators being

considered, it is a simple matter to numerically minimize over  $M$  to find the optimal  $M^*$ . Table 1 presents the values of  $M^*$  for various  $\epsilon$  and  $p = 3, 5, 7$ . Thus, when  $\epsilon = .2$ ,  $p = 3$ , and  $\Gamma_2$  is the class of priors considered, the estimator of the form (2.5) which minimizes  $r_{\Gamma_2}^*$  is  $\delta_{\delta}^{.42, S}$  (i.e.,  $M^* = .42$ ).

Table 2 compares the robustness of  $\delta_{\delta}^{M^*, I}$  and  $\delta_{\delta}^{M^*, S}$ , as measured by  $r_{\Gamma_1}^*$ , for various  $\epsilon$  and  $p = 3, 5$ , and  $7$ . (The estimator  $\delta_{\delta}^T$  will be discussed shortly.) Table 3 presents the analogous results for robustness as measured by  $r_{\Gamma_2}^*$ . In both cases,  $\delta_{\delta}^{M^*, I}$  is superior to  $\delta_{\delta}^{M^*, S}$ . Thus greater overall robustness appears to be available from the class of coordinatewise independent estimators (i.e. (1.7)), than from the class of Stein-effect estimators. This is in sharp contrast to the results obtained in Berger (1982a, b, 1984) for the "dependent" class  $\Gamma$  in (1.4).

Of course,  $r_{\Gamma_1}^*$  and  $r_{\Gamma_2}^*$  do not tell the whole story, in that they essentially measure only the "worst" that can happen. It is thus interesting to also look at the "best" that can happen in using an estimator  $\delta$ ; this best will be roughly given by  $RSR(\pi^0, \delta)$ , which again can be thought of as the percentage of available Bayesian gains sacrificed by using  $\delta$  instead of  $\delta_{\delta}^{\pi^0}$ . Table 4 gives values of RSR for  $\delta_{\delta}^{M^*, I}$  and  $\delta_{\delta}^{M^*, S}$ , when  $M^*$  is the optimal choice with respect to  $\Gamma_1$ , while Table 5 gives the corresponding values of RSR, when  $M^*$  is the optimal choice with respect to  $\Gamma_2$ . Observe that the Stein-effect estimators seem to do better than the coordinatewise independent estimators. Thus, while  $\delta_{\delta}^{M^*, I}$  is generally better than  $\delta_{\delta}^{M^*, S}$  in the "worst case scenario," the reverse seems to be true in the "best case."

To attempt a finer comparison of the coordinatewise independent and Stein-effect estimators in terms of RSR, it was decided to replace  $M^*$  in  $\delta_{\delta}^{M^*, I}$  by  $M'$  the value of  $M$  for which  $\delta_{\delta}^{M', I}$  and  $\delta_{\delta}^{M', S}$  are equally robust (i.e., for which  $r_{\Gamma}^*(\delta_{\delta}^{M', I}) = r_{\Gamma}^*(\delta_{\delta}^{M', S})$ ). The RSR of  $\delta_{\delta}^{M', I}$  was then calculated, and is given in



Tables 4 and 5. The results are unclear. For small  $\varepsilon$  the coordinatewise independent estimator,  $\hat{\delta}^{M^I}$ , tends to be better, while for large  $\varepsilon$  the Stein-effect estimator,  $\hat{\delta}^{M^S}$ , has smaller RSR. This is, perhaps, not surprising, in that, at the extreme  $\varepsilon = 1$ , only minimax estimators are allowed, and only Stein-effect estimators can be minimax and have small RSR.

A different method for achieving a type of Bayesian robustness was proposed in Stein (1981) (see also Dey and Berger (1983)). The method, when applied to our situation, results in the estimator

$$\hat{\delta}^T(x) = x - \min \left\{ \frac{\sigma^2}{(\sigma^2 + \tau^2)}, \frac{2(\ell - 2)}{|z|_2^2} \right\} z;$$

here  $z = (z_1, \dots, z_p)^t$ ,  $z_i = \text{sgn}(y_i) \min\{|y_i|, |y|_{(\ell)}\}$ ,  $y_i = (x_i - \mu_i)$ , and  $|y|_{(\ell)}$  is the  $\ell^{\text{th}}$  order statistic of the sequence  $\{|y_1|, |y_2|, \dots, |y_p|\}$ .

The choice of  $\ell$  should correspond roughly to a guess as to the number of  $\theta_i$  for which prior misspecification is feared; because of the small  $p$  in this study and because  $\Gamma_2$  corresponds to the case where at most one coordinate is misspecified, the choice  $\ell = p-1$  was made.

Using results of Dey and Berger (1983), RSR  $(\pi^0, \hat{\delta}^T)$  can be calculated.

For  $p = 5$  and  $p = 7$ , the values are .383 and .337, respectively. Furthermore

$r_{\Gamma_1}^*(\hat{\delta}^T)$  and  $r_{\Gamma_2}^*(\hat{\delta}^T)$  can be calculated in a manner similar to (but more complicated than) Lemmas 1 and 2. The key idea is to recognize that, if  $q_i$  gives

probability one to  $\theta_i^*$  and  $|\theta_i^*| \rightarrow \infty$ , then  $\hat{\delta}^T$  effectively becomes the estimator with  $z_j = y_j$  for  $j \neq i$ , and  $z_i = \text{sgn}(y_i) \max_{j \neq i} |y_j|$ ; although this is the worst

case, a Stein effect still remains. If, on the other hand, one is maximizing over two contaminations  $q_i$  and  $q_j$ , they can be chosen to give probability one to

$\theta_j^*$  and  $\theta_j^*$ , respectively, with  $|\theta_j^*| \rightarrow \infty$  and  $|\theta_j^*| \rightarrow \infty$ ;  $\delta_\epsilon^T$  will then collapse back to  $\delta_\epsilon^0(\underline{x}) = \underline{x}$ . Thus, in the expansion of terms in  $r_{\Gamma_1}^*$  or  $r_{\Gamma_2}^*$ , those involving at most one  $q_j$  yield reduced Bayes risk, while those involving two or more of the  $q_j$  result in the minimax risk. Expressions in Dey and Berger (1983) can be used to calculate all terms.

Tables 2 and 3 give values of  $r_{\Gamma_1}^*$  and  $r_{\Gamma_2}^*$  for  $\delta_\epsilon^T$ , when  $p = 5$  or  $7$ . (For  $p = 3$ ,  $\delta_\epsilon^T(\underline{x}) = \delta_\epsilon^0(\underline{x}) = \underline{x}$ , since  $\ell = p - 1 = 2$ .) Even for  $\Gamma_2$  (a class of priors for which  $\delta_\epsilon^T$  would seem ideally suited, since  $\delta_\epsilon^T$  truncates precisely one "bad" coordinate here), the performance of  $\delta_\epsilon^T$  seems inferior to that of the robust coordinatewise independent estimator. Its value of  $r_{\Gamma}^*$  is always larger, and its RSR (.383 or .337 for  $p = 5$  or  $7$ ) is better only for quite large  $\epsilon$ . Observe that  $\delta_\epsilon^T$  does seem to be better than  $\delta_\epsilon^{M^*,S}$  in terms of  $r_{\Gamma}^*$ , unless  $\epsilon$  is small, but its RSR is never better than that of  $\delta_\epsilon^{M^*,S}$ .

#### 4. CONCLUSIONS

The numerical results in Section 3 indicate that a Bayesian who is quite certain that several problems are unrelated, but is otherwise somewhat uncertain about the prior specification, will not see any clear value in utilizing the Stein effect; superior robustness can be achieved through use of, say, coordinatewise independent limited translation estimators (except possibly for large  $\epsilon$ , where no clear conclusions emerged from the limited criteria considered). The coordinatewise independent estimators also have several additional appealing properties. First, in contrast with the Stein-effect estimators (though not necessarily the  $\delta_\epsilon^T$  versions), the coordinatewise independent estimators do not suffer from the problem that one or more badly specified  $\pi_i^0$  can cause the entire estimator to collapse back to  $\underline{x}$ . (A single large  $|x_i - \mu_i|$  will cause  $\delta_\epsilon^{M,S}$  to

be approximately  $\chi$ ; one bad apple can thus ruin the pie.) A related type of "additional robustness," of the coordinatewise independent estimators, is robustness with respect to the loss. The coordinatewise risk of  $\delta_{\chi}^{M^*,S}$  can be quite large (cf. Efron and Morris (1973)), while the coordinatewise risk of  $\delta_{\chi}^{M^*,I}$  is inherently controlled; thus, if there is uncertainty about the "weights" that should be assigned to each component of the loss,  $\delta_{\chi}^{M^*,I}$  can be substantially more robust. ( $\delta_{\chi}^T$  will be substantially more robust than  $\delta_{\chi}^{M^*,S}$  in terms of this type of loss robustness, but, again,  $\delta_{\chi}^T$  sacrifices a good deal in terms of  $r_{\Gamma}^*$  and RSR.) Related to this loss robustness is the fact that there is no need to determine weights for the coordinatewise losses,  $L_i(\theta_i, \delta_i)$ , in calculating  $\delta_{\chi}^{M^*,I}$  (just do each coordinate separately), while the optimal Stein-effect estimator will depend crucially on such weights. From a practical perspective this might be of paramount concern, since it is often very difficult to convince practitioners to compare and select weights for the losses from the separate coordinate problems.

Note that the above conclusions have some force for non-Bayesians also, because of the nature of Stein estimation. It can be argued (see Berger and Berliner (1984) and the discussion by Brown in Berger (1984)) that (i) prior information must be utilized to select an alternative to  $\delta_{\chi}^0$ ; (ii) the most appealing method, in Stein estimation, of utilizing such information is to specify a prior distribution,  $\pi^0$ , and seek to minimize  $r(\pi^0, \delta_{\chi})$  subject to some frequentist risk restriction (say,  $R(\theta_{\chi}, \delta_{\chi}) \leq C$ ); and (iii) a more intuitively accessible mathematical formulation of this last problem is the  $\Gamma$ -minimax problem of selecting a  $\Gamma$  as in (1.4), (1.8), or (2.1) (with  $\varepsilon$  related to  $C$ ), and finding  $\delta_{\chi}$  to minimize  $r_{\Gamma}^*(\delta_{\chi})$ . Thus following a natural path within frequentist Stein-estimation theory leads to consideration of the issues we have been discussing.

Of course, the conclusion that the Stein effect is of no clear value in

combining completely unrelated problems does not say that it is of no use in combining related problems. And, from a practical perspective, the only problems in which practitioners would even think of using a combined estimator are those in which the parameters might well be related. The message, therefore, is simply that there is no "paradoxical" feature of Stein estimation to a Bayesian. If there is reason to suspect some prior relationship between the  $\theta_i$ , a Bayesian might well use a (suitably tailored) Stein-effect estimator; otherwise he will not. The choice would involve careful consideration of prior knowledge.

Even when the  $\theta_i$  are clearly independent a priori, a Bayesian consultant can utilize the Stein-effect when dealing with frequentist clients. The use of  $\delta^{M^*,S}$  can be shown to be preferable to the use of  $\delta^0(x) = x$  from a Bayesian perspective (even conditionally);  $\delta^{M^*,S}$  can thus be recommended when frequentist justification is an externally imposed requirement. Of course, a Bayesian might feel somewhat silly in recommending the combining of completely unrelated problems in such a scenario, but that is the price that would have to be paid for agreeing to guarantee standard frequentist validity.

#### BIBLIOGRAPHY

- Berger, J., (1982a). Bayesian robustness and the Stein effect. J. Amer. Statist. Assoc., 77, 358-368.
- Berger, J., (1982b). Estimation in continuous exponential families: Bayesian estimation subject to risk restrictions and inadmissibility results. In Statistical Decision Theory and Related Topics III, S. S. Gupta and J. Berger (Eds.). Academic Press, New York.
- Berger, J., (1984). The robust Bayesian viewpoint. In Robustness in Bayesian Statistics, J. Kadane (Ed.). North-Holland, Amsterdam.
- Berger, J., (1985a). The Stein effect. In Encyclopedia of Statistical Sciences, S. Kotz and N. L. Johnson (Eds.). Wiley, New York.
- Berger, J., (1985b). Statistical Decision Theory and Bayesian Analysis. Springer-Verlag, New York.

- Berger, J., and Berliner, L. M., (1984). Bayesian input in Stein-estimation and a new minimax empirical Bayes estimator. J. of Econometrics, 25, 87-108.
- Chen, S. Y., (1983). Restricted risk Bayes estimation. Ph.D. Thesis, Department of Statistics, Purdue University, West Lafayette.
- Dey, D. K., and Berger, J., (1983). On truncation of shrinkage estimators in simultaneous estimation of normal means. J. Amer. Statist. Assoc., 78, 865-869.
- Efron, B., and Morris, C., (1971). Limiting the risk of Bayes and empirical Bayes estimators--Part I: the Bayes case. J. Amer. Statist. Assoc., 66, 807-815.
- Efron, B., and Morris, C., (1973). Stein's estimation rule and its competitors--an empirical Bayes approach. J. Amer. Statist. Assoc., 68, 117-130.
- James, W., and Stein, C., (1960). Estimation with quadratic loss. In Proc. Fourth Berkeley Symposium on Math. Statist. and Prob., 1, 361-379. University of California Press, Berkeley.
- Marazzi, A., (1985). On constrained minimization of the Bayes risk for the linear model. Statistics and Decisions.
- Stein, C., (1955). Inadmissibility of the usual estimator for the mean of a multivariate normal distribution. In Proc. Third Berkeley Symp. Math. Statist. Prob., 1, 197-206. University of California Press, Berkeley.
- Stein, C., (1981). Estimation of the mean of a multivariate normal distribution. Ann. Statist., 9, 1135-1151.

Table 1. Values of  $M^*$ .

$\frac{p}{\epsilon}$	$\Gamma_{1,\delta}^{M^*,I}$	$\Gamma_{2,\delta}^{M^*,I}$	$\Gamma_{1,\delta}^{M^*,S}$			$\Gamma_{2,\delta}^{M^*,S}$		
	all p	all p	3	5	7	3	5	7
0.0	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$	$\infty$
0.1	1.30	2.38	0.17	0.002	0*	0.50	0.19	0.05
0.2	0.74	1.67	0.05	0*	0*	0.24	0.06	0*
0.3	0.47	1.30	0.02	0*	0*	0.14	0.02	0*
0.4	0.30	1.06	0.007	0*	0*	0.09	0.002	0*
0.5	0.19	0.88	0.001	0*	0*	0.05	0*	0*
0.6	0.11	0.74	0.001	0*	0*	0.03	0*	0*
0.7	0.06	0.63	0*	0*	0*	0.01	0*	0*
0.8	0.03	0.54	0*	0*	0*	0.006	0*	0*
0.9	0.006	0.47	0*	0*	0*	0.001	0*	0*
1.0	0	0	0	0	0	0	0	0

0\* means smaller than  $1 \times 10^{-4}$ .







Table 4. Values of  $RSR(\pi^0, \hat{\rho})$  for  $\hat{\rho}^{M^*, I}$ ,  $\hat{\rho}^{M^*, S}$ , and  $\hat{\rho}^{M^I, I}$ , when  $M^*$  is optimal for  $\Gamma_1$ .

$\epsilon \backslash p$	$RSR(\pi^0, \hat{\rho}^{M^*, I})$	$RSR(\pi^0, \hat{\rho}^{M^*, S})$			$RSR(\pi^0, \hat{\rho}^{M^I, I})$		
	all p	3	5	7	3	5	7
0.0	0	0	0	0	0	0	0
0.1	0.11	0.10	0.07	0.03	0.03	0.02	0.01
0.2	0.20	0.11	0.07	0.03	0.08	0.04	0.03
0.3	0.29	0.21	0.07	0.03	0.10	0.08	0.06
0.4	0.38	0.24	0.07	0.03	0.18	0.13	0.11
0.5	0.47	0.27	0.07	0.03	0.26	0.20	0.19
0.6	0.57	0.28	0.07	0.03	0.35	0.30	0.29
0.7	0.66	0.29	0.07	0.03	0.46	0.42	0.42
0.8	0.77	0.29	0.07	0.03	0.60	0.58	0.58
0.9	0.88	0.29	0.07	0.03	0.77	1.00	0.77
1.0	1.00	0.29	0.07	0.03	1.00	1.00	1.00

Table 5. Values of  $RSR(\pi^0, \hat{\rho})$  for  $\hat{\rho}^{M^*, I}$ ,  $\hat{\rho}^{M^*, S}$ , and  $\hat{\rho}^{M', S}$ , when  $M^*$  is optimal for  $\Gamma_2$ .

$\frac{p}{\epsilon}$	$RSR(\pi^0, \hat{\rho}^{M^*, I})$			$RSR(\pi^0, \hat{\rho}^{M^*, S})$			$RSR(\pi^0, \hat{\rho}^{M', S})$		
	3	5	7	3	5	7	3	5	7
0.0	0	0	0	0	0	0	0	0	0
0.1	0.04	0.03	0.02	0.04	0.03	0.02	.009	.003	.001
0.2	0.08	0.05	0.04	0.08	0.05	0.03	.018	.005	.002
0.3	0.11	0.07	0.06	0.11	0.06	0.03	.025	.007	.002
0.4	0.14	0.09	0.08	0.14	0.07	0.03	.032	.008	.002
0.5	0.17	0.11	0.09	0.17	0.07	0.03	.038	.009	.003
0.6	0.20	0.13	0.10	0.20	0.07	0.03	.044	.010	.003
0.7	0.23	0.15	0.11	0.22	0.07	0.03	.049	.011	.003
0.8	0.26	0.17	0.12	0.25	0.07	0.03	.054	.011	.003
0.9	0.29	0.19	0.14	0.27	0.07	0.03	.059	.012	.003
1.0	1.00	1.00	1.00	0.29	0.07	0.03	1.00	1.00	1.00