

Loss Functions for Set Estimations

BU-999-MA

June, 1992

George Casella, J.T.-Gene Hwang<sup>1</sup>, and Christian P. Robert<sup>2</sup>

Cornell University, Ithaca, N.Y. and Paris University (Rouen University?)

Key Words and Phrases: Confidence sets, Decision theory, Bayes estimation.

AMS 1980 Subject Classification: Primary 62C05, Secondary 62F25, 62A99.

---

<sup>1</sup> Research supported by National Science Foundation Grant No. DMS88-09016.

<sup>2</sup> Research supported by the U.S. Army Research Office through the Mathematical Sciences Institute at Cornell University.

### Summary

Approaches to set estimation based on a decision-theoretic formulation have usually used a loss function that is a linear combination of volume and coverage probability. Such loss functions can suffer from paradoxical behavior of the Bayes rules, and thus may not be appropriate. We investigate the behavior of optimal set estimators for different classes of loss functions and study their decision-theoretic properties.

## 1. Introduction

Set estimators contribute an inherent part of most statistical analyses, providing sets  $C_x$  depending on the observation  $x \sim f(x|\theta)$  for the parameter of interval  $\theta$ . Although they are thus bound to be used extensively in decisions, there has been hardly any move towards a decision-theoretic evaluation of these procedures.

Set estimators are usually derive from testing acceptance regions by duality. However, contrary to intuition, these estimators do not necessarily enjoy the optimality properties of their testing dual sets, as shown by Hwang and Casella (1982) for the domination of the usual normal confidence region.

A direct approach to set estimation has to take into account two criteria, size and coverage. The resting alternative can also be perceived this way for one of the two criteria remaining constant, namely the coverage probability of the considered set. Apart from vector loss generalizations, for which a complete ordering is most often impossible, a loss function considered in the literature is the *linear loss function*,

$$(1.1) \quad L(\theta, C) = a \text{ vol}(C) + \mathbb{I}(\theta \in C) ,$$

where  $a$  is a fixed constant, as in Cohen and Strawderman (1973) and Meeden and Vardeman (1985). This loss function indeed allows for interplay between volume and coverage but it can also suffer from major defects. One of these leads to a paradox pointed to us by James Berger and studied in Casella, Hwang and Robert (1992). We also refer the reader to this paper for additional references on loss estimation.

The so-called paradox exhibits a domination under (1.1) of the Student's t-interval by a truncated set which is empty when  $s^2$ , the empirical variance, is *large enough*. However, the Bayes estimator associated with (1.1) and the Lebesgue prior is such that its length is decreasing in  $S$  when  $S$  is large enough. These phenomena incriminate the loss (1.1) and the fact that it does not balance properly the two criteria, giving an undue importance to the volume component.

In this paper, we propose to study the decision-theoretic properties of the losses of the form

$$(1.2) \quad L(\theta, C) = S[\text{vol}(C)] + \Pi(\theta \notin C) ,$$

where  $S$  is a *size* function. Casella *et al.* (1992) already showed that some losses in (1.2) can avoid Berger's paradox and came to the following requirements on  $S$ : the size function must be increasing and satisfy  $S(0) = 0$ ,  $\lim_{b \rightarrow +\infty} S(b) = 1$  so that volume *and* coverage are weighted equally. These conditions are only necessary since some functions  $S$  satisfy them but still lead to Berger's paradox (see Casella *et al.*, 1992).

Section 2 is dedicated to a classical analysis of losses like (1.2) and the derived notions of admissibility, minimaxity and Bayes data. Section 3 examines the range of Bayes sets for different parameterized families of losses and derives conditions under which conjugate Bayes sets are never empty. Section 4 focuses on two particular classes of losses, one the bounded spaces, the other for unbounded spaces.

## 2. Decision Theory

In this section, we consider usual decision-theoretic properties derived from the loss function

$$(2.1) \quad L_{\mathfrak{S}}(\theta, C) = \mathfrak{S}[\text{vol}(C)] + \mathbb{I}(\theta \notin C) ,$$

where  $\mathfrak{S}$  is increasing from 0 to 1. First note that, as mentioned in Joshi (1969), the action space is not  $\mathfrak{B}(\Theta)$ , set of all measurable  $v$  subsets of  $\Theta$ , but rather  $\mathbb{B}(\Theta)$ , the quotient space derived from  $\mathfrak{B}(\Theta)$  by the equivalence relation  $C \sim C'$  iff  $\lambda(C \Delta C') = 0$ , where  $\lambda$  is the Lebesgue measure on  $\Theta$  and  $C \Delta C'$  denotes  $(C \setminus C') \cup (C' \setminus C)$ . In other words, confidence sets can only be defined up to a set of Lebesgue measure zero. Otherwise, no comparison is possible. For this reason, the only acceptable priors are those which are absolutely continuous with the Lebesgue measure.

### 2.1. Bayes and Generalized Bayes Estimators

For any given loss function  $L(\theta, C)$  and prior distribution  $\pi$  on  $\theta$ , a (decision-theoretic) *Bayes estimator*  $C^\pi$  is a solution of the minimization problem

$$(2.2) \quad \int_{\theta} R(\theta, C^\pi) \pi(\theta) d\theta = \min_C \int_{\theta} R(\theta, C) \pi(\theta) d\theta .$$

If  $\pi(\theta)$  is a proper prior distribution, for each  $x$  the Bayes estimator  $C_x^\pi$  minimizes

$$(2.3) \quad \int_{\theta} L(\theta, C) \pi(\theta | x) d\theta ,$$

where  $\pi(\theta | x)$  is the posterior distribution. If  $\pi(\theta)$  is not a proper prior distribution but  $\pi(\theta | x)$  is proper, then (2.3) is taken as the definition of  $C_x^\pi$ . Thus, in general,  $C_x^\pi$  minimizes the posterior expected loss. For loss functions (2.1), the Bayes estimators are HPD regions.

**Theorem 2.1.** *For  $X \sim f(x | \theta)$  and (possibly improper) prior  $\pi(\theta)$ , the Bayes estimator against the loss (2.1) is given by*

$$C_x^\pi = \left\{ \theta : \pi(\theta) \geq k(x) \right\} ,$$

where  $k = k(x)$  minimizes the quantity

$$\mathfrak{S}\left(\text{vol}\{\theta : \pi(\theta | x) \geq k\}\right) - \int_{\{\theta : \pi(\theta | x) \geq k\}} \pi(\theta | x) d\theta ,$$

provided that  $\left\{ \theta : \pi(\theta | x) = \ell \right\}$  has Lebesgue measure zero for every  $\ell$ .

The condition that  $\{\theta: \pi(\theta | x) = \ell\}$  has zero Lebesgue measure is a technical one, and has been added to eliminate the necessity of considering randomized rules. We chose to include it in the assumptions for ease of understanding. A more general theorem which does not need this assumption is included, for completeness, in the Appendix. In most cases, however, Theorem 2.1 is general enough.

## 2.2 Admissible Estimators

Given the previous restrictions on the definition of the action space, we will say that an estimator  $C$  is *admissible* if, whenever there exists an estimator  $C'$  satisfying  $R(\theta, C') \leq R(\theta, C)$  for almost all  $\theta$ , then  $R(\theta, C') = R(\theta, C)$  for almost all  $\theta$ .

We will consider here the relation between Bayes estimators and admissibility. Things are not as clear as for point estimation and some questions remain unanswered; in particular, we do not yet know if there exist equivalent theorems to the complete class theorems of point estimation.

**Theorem 2.2.** *If  $\pi$  gives positive measure to any set with positive Lebesgue measure, then the Bayes estimator  $C^\pi$  is admissible.*

If  $\pi$  satisfies the conditions of the previous theorem, then all the Bayes sets are admissible, even if  $C^\pi$  is not essentially unique. However, unlike the point estimation case with strictly convex loss, it may be possible to exhibit inadmissible proper Bayes estimators. Obviously, generalized Bayes estimators may also be inadmissible (as in the point estimation problem).

**Example 2.1.** It is easily seen that the usual confidence set

$$C_x^0 = \{\theta: |\theta - x| \leq c\}$$

is generalized Bayes with respect to the Lebesgue measure. However, for any choice of the loss,  $C_x^0$  is inadmissible. Consider the recentered set,

$$C_x^a = \{\theta: |\theta - \delta_a^+(x)| \leq c\},$$

where  $\delta_a^+$  is the positive-part James Stein estimator

$$\delta_a^+(x) = \left(1 - \frac{a}{|x|^2}\right)^+ x.$$

Hwang and Casella (1982, 1984) have established that  $P_\theta(\theta \in C_x^0) < P_\theta(\theta \in C_x^a)$  for all  $\theta$  and some values of  $a$  (depending on  $c$ ), and the two sets  $C_x^0$  and  $C_x^a$  have the same volume.

It is not true that any  $S(\cdot)$  function considered in (2.1) will result in nonparadoxical behavior. There are, in fact, numerous examples of losses and prior distributions which give  $\phi$  (or  $\Theta$ ) as an admissible constant Bayes set estimator. This problem is considered in Section 4.

### 2.3. Minimavity

For a given loss,  $L$ , the *minimax risk* is defined by

$$(2.4) \quad \mathbf{m} = \inf_C \sup_{\theta} R(\theta, C).$$

As  $R(\theta, C) = 1$  for every  $\theta$ , it is obvious that  $\mathbf{m} \leq 1$ . We would like to know when  $\theta$  is a minimax estimator, i.e., when does  $\mathbf{m} = 1$ , in order to avoid such losses. In particular, we have the following condition.

**Theorem 2.3.** *If  $X \sim N_p(\theta, I)$ , a necessary and sufficient condition for  $\mathbf{m} = 1$  is that*

$$S(\mathfrak{B}_{c^p}) \geq P(|Z| < c)$$

for every  $c \geq 0$ , where  $Z \sim N_p(0, I)$  and  $\beta$  is the volume of the unit ball.

**Proof.** If there exists  $c_0$  such that  $S(\mathfrak{B}_{c_0^p}) < P(|Z| < c_0)$ , the estimator

$$C_x^0 = \{\theta : |x - \theta| \leq c_0\}$$

has a constant risk strictly smaller than 1. Hence  $\mathbf{m} < 1$ , establishing the necessity. For the sufficiency, consider the conjugate priors  $N(0, \tau^2 I)$ . The corresponding Bayes sets are

$$C_x^\pi = \left\{ \theta : \left| \theta - \frac{\tau^2}{\tau^2 + 1} x \right| < \ell_\tau^* \sqrt{\frac{\tau^2}{\tau^2 + 1}} \right\}$$

where  $\ell_\tau^*$  attains the minimum

$$\min_{\ell} \left[ S \left( \mathfrak{B} \left( \frac{\tau^2}{\tau^2 + 1} \right)^{p/2} \ell^p \right) + P(|Z| > \ell) \right].$$

The minimum value is also the Bayes risk and hence is a lower bound on  $\mathbf{m}$ . The sufficiency will therefore be established if we can show that the minimum is bounded below by zero as  $\tau^2 \rightarrow \infty$ . Now the minimum equals

$$\begin{aligned} \min_{\ell} \left[ \mathbf{S}(\mathfrak{B}\ell^p) - \mathbb{P} \left( |Z| < \left( \frac{\tau^2}{\tau^2 + 1} \right)^{-1/2} \ell \right) \right] &\geq \min_{\ell} \left[ \mathbf{S}(\mathfrak{B}\ell^p) - \mathbb{P}(|Z| < \ell) \right] \\ &+ \min_{\ell} \left[ \mathbb{P}(|Z| < \ell) - \mathbb{P} \left( |Z| < \left( \frac{\tau^2 + 1}{\tau^2} \right)^{1/2} \ell \right) \right]. \end{aligned}$$

The first term of the lower bound is nonnegative by assumption, and the second term goes to zero as  $\tau^2 \rightarrow \infty$ . The proof is now complete.  $\square$

An example of loss for which  $\mathbf{m} = 0$  is given below. First, we define the effective radius of a  $p$  dimensional set  $C$ ,  $\rho(C)$ , as  $\rho(C) = (\text{vol}(C)/\mathfrak{B})^{1/p}$ .

**Example 2.2.** Let  $X \sim \mathcal{N}(\theta, I_p)$ ,  $Z \sim \mathcal{N}(0, I_p)$  and the loss is

$$L(\theta, C) = \mathbb{P}(|Z| < 2\rho(C)) + \mathbb{I}(\theta \notin C),$$

under the prior distribution  $\mathcal{N}(0, \tau^2)$ . Then the posterior distribution is  $\mathcal{N}\left(\frac{\tau^2}{\tau^2 + 1}x, \frac{\tau^2}{\tau^2 + 1}I\right)$  and the Bayes set is of the form

$$C_x^\pi = \left\{ \theta : \left| \theta - \frac{\tau^2}{\tau^2 + 1}x \right| \leq \ell^* \sqrt{\frac{\tau^2}{\tau^2 + 1}} \right\},$$

where  $\ell^*$  minimizes

$$(2.5) \quad \mathbb{P} \left( |Z| \leq 2\ell \sqrt{\frac{\tau^2}{\tau^2 + 1}} \right) - \mathbb{P}(|Z| \leq \ell).$$

If  $[\tau^2/(\tau^2 + 1)] > 1/4$ , a solution of (2.5) is  $\ell^* = 0$  or  $\infty$ ; thus  $C_x^\pi = \emptyset$  or  $\Theta$ . Since  $\pi$  satisfies the conditions of Theorem 2.2,  $\emptyset$  is admissible and  $\mathbf{m} = 0$ .

The above example illustrates a technique used to establish that  $\mathbf{m} = 1$ , namely to show that some Bayes estimators are  $\phi$ . The nonexistence of trivial Bayes sets is actually a necessary condition for  $\mathbf{m} < 1$  (see Section 2). Another way to show that  $\mathbf{m} < 1$  is to exhibit an estimator with maximum risk (in  $\theta$ ) less than 1, as illustrated in the following example, for a bounded parameter space  $[0, 1]$ . In such a case, size function is the length (volume).

**Example 2.3.** Let  $X \sim \text{Binomial}(n, \theta)$ . Then, because  $\theta \in [0, 1]$ , the simple loss

$$(2.6) \quad L(\theta, C) = \ell(C) + \mathbb{I}(\theta \notin C),$$

where  $\ell$  is the length of  $C$ , fits the requirements of (2.1). If  $\pi(\theta)$  is the uniform distribution on  $[0, 1]$ ,

then  $\pi(\theta | \mathbf{x})$  is Beta( $x + 1, n - x + 1$ ). We will see in Section 5 that the Bayes rule is necessarily the set  $\{\pi(\theta | \mathbf{x}) \geq 1\}$ , which is an interval  $[\theta_{\mathbf{x}}^{\ell}, \theta_{\mathbf{x}}^{\text{u}}]$ . Numerical solutions show that, for all the observed values of  $n$ , the frequentist risk is strictly negative for every  $\theta \in [0, 1]$ . Therefore, under (2.6),  $\mathbf{m} < 1$ . Note that the Bayes sets  $[\theta_{\mathbf{x}}^{\ell}, \theta_{\mathbf{x}}^{\text{u}}]$  are not the UMPU intervals (Blyth and Hutchinson, 1961). For selected values of  $n$ , Table 2.1 gives the minimum coverage probability and maximum length of the Bayes rule.

Table 2.1. Minimum coverage probability and maximal length of the Bayes set for binomial  $\theta$  using a uniform prior.

n	10	20	25	30	35	40	45
Length	.41	.34	.32	.30	.29	.28	.27
Posterior coverage	.79	.87	.88	.90	.91	.91	.92

However, the Bayes rule is not minimax as the associated Bayes risk is always strictly lower than the frequentist risk.

### 3. Ranges of Bayes Sets

As previously seen, one problem with a loss function approach to set estimation is the range of resulting Bayes sets. In particular, it is possible to derive proper Bayes sets that are trivial, that is, Bayes sets that may be either empty or equal to the entire parameter space. In Section 3.1, we derive a sufficient condition for a Bayes set to be nontrivial. We say that the Bayes set is *nontrivial* if it is neither  $\emptyset$  nor  $\Theta$  with positive posterior probability. In Section 3.2 we consider particular classes of losses and study the ranges of the Bayes sets. We give explicit bounds on this range in the normal case.

#### 3.1. Conditions for Existence of Nontrivial Bayes Sets

We can derive a useful sufficient condition by differentiating

$$(3.1) \quad \mathbf{S}\left(\text{vol}\left(C_x^\pi(k)\right)\right) - \int_{C_x^\pi(k)} \pi(\theta|x)d\theta ,$$

where  $C_x^\pi(k) = \{\theta : \pi(\theta|x) > k\}$ , with respect to  $k$ .

**Lemma 3.1.** *For the set  $C_x^\pi(k) = \{\theta : \pi(\theta|x) \geq k\}$ , is  $\mathbf{S}$  is differentiable, then*

$$(3.2) \quad \frac{\partial}{\partial k} \left[ \mathbf{S}\left(\text{vol}\left(C_x^\pi(k)\right)\right) - \int_{C_x^\pi(k)} \pi(\theta|x)d\theta \right] = \left[ k - \mathbf{S}'\left(\text{vol}\left(C_x^\pi(k)\right)\right) \right] \int_{\{\theta:\pi(\theta|x)=k\}} \frac{ds}{|\nabla\pi(\theta|x)|} .$$

Here,  $ds$  represents the infinitesimal surface area of the set  $\{\theta : \pi(\theta|x) = k\}$  and  $\nabla\pi(\theta|x)$  is the gradient of  $\pi(\theta|x)$  for fixed  $x$ .

**Proof.** Straightforward differentiation yields

$$\frac{\partial}{\partial k} \left[ \mathbf{S}\left(\text{vol}\left(C_x^\pi(k)\right)\right) - \int_{C_x^\pi(k)} \pi(\theta|x)d\theta \right] = \mathbf{S}'\left(\text{vol}\left(C_x^\pi(k)\right)\right) \frac{\partial}{\partial k} \text{vol}\left(C_x^\pi(k)\right) - \frac{\partial}{\partial k} \int_{C_x^\pi(k)} \pi(\theta|x)d\theta .$$

From the definition of  $C_x^\pi(k)$ ,

$$(3.2) \quad \frac{\partial}{\partial k} \text{vol}\left(C_x^\pi(k)\right) = \frac{\partial}{\partial k} \int_{\{\theta:\pi(\theta|x) \geq k\}} d\theta = \lim_{\Delta k \rightarrow 0} \frac{1}{\Delta k} \int_{\{\theta:k+\Delta k \geq \pi(\theta|x) \geq k\}} d\theta .$$

Let  $h_x(\theta)$  be the perpendicular infinitesimal distance between the sets  $\{\theta : \pi(\theta|x) = k\}$  and  $\{\theta : \pi(\theta|x) = k + \Delta k\}$ . Clearly  $\Delta k = h_x(\theta) |\nabla\pi(\theta|x)|$  and hence from (3.2) we obtain

$$(3.4) \quad \frac{\partial}{\partial k} \text{vol}(C_X^\pi(k)) = \lim_{\Delta k \rightarrow 0} \frac{1}{\Delta k} \int h_x(\theta) ds = - \int \frac{ds}{|\nabla \pi(\theta|x)|}.$$

Similarly, we have

$$(3.5) \quad \frac{\partial}{\partial k} \int_{\{\theta: \pi(\theta|x) \geq k\}} \pi(\theta|x) d\theta = -k \int \frac{ds}{|\nabla \pi(\theta|x)|}$$

which, together with (3.4), establishes the lemma.  $\square$

One nice feature of equation (3.2) is that to evaluate the sign of this derivative, only the term  $[k - S'(\text{vol}(C_X^\pi(k)))]$  is important, as the sign of the integral is always nonnegative. Therefore we only need this function to determine whether an interior point is a minimum or maximum and the difficult task of calculating this surface integral can be avoided.

A useful sufficient condition can now be easily established. Due to the fact that the posterior loss (3.1) is equal to zero for  $k = 0$  and  $k = \infty$ , straightforward application of Lemma 2.1 yields the following sufficient conditions.

**Theorem 3.2.** *If there exists  $k_0 > 0$  such that*

$$[k - S'(\text{vol}(C_X^\pi(k)))] < 0 \quad \text{for } k \text{ satisfying } 0 < k < k_0,$$

or

$$[k - S'(\text{vol}(C_X^\pi(k)))] > 0 \quad \text{for all } k > k_0,$$

the Bayes set  $C_X^\pi$  is nontrivial.

### 3.2. The Normal Case

We now investigate estimation of  $\theta$  in the normal setting with

$$X|\theta \sim N_p(\theta, \sigma^2 I), \quad \theta \sim N_p(\mu, \tau^2 I), \quad 0 < \tau^2 \leq \infty.$$

The posterior distribution of  $\theta|X = x$  is  $N_p(m(x), \sigma_\tau^2 I)$ , with

$$m(x) = \frac{\sigma_\tau^2}{\tau^2} \mu + \left(1 - \frac{\sigma_\tau^2}{\tau^2}\right)x, \quad \sigma_\tau^2 = \frac{\sigma^2 \tau^2}{\sigma^2 + \tau^2}.$$

Instead of using Theorem 3.2, we take a slightly different approach in this setting. We feel that this variation, which is only appropriate in the normal setting, will provide more insight into the problem.

Recall the definition of the effective radius  $\rho(C)$ , of a set  $C$ . We consider three classes of size functions expressed in terms of effective radius  $\rho = \rho(C)$ :

$$\begin{aligned} \text{i) } \bar{S}_a(\rho) &= \rho^p / (a + \rho^p), & a > 0 \\ \text{ii) } \bar{S}_a(\rho) &= P(|Z| \leq a\rho), & a > 0, Z \sim N_p(0, 1) \\ \text{iii) } \bar{S}_a(\rho) &= 1 - e^{-a\rho^2/2}, & a > 0. \end{aligned}$$

It appears that the class of losses defined using (3.6i) is a reasonable class for most situations (see Section 4). Examination of the losses given by (3.6ii) and (3.6iii) will further reveal the bounds of the theory considered here.

We can write HPD regions for  $\theta$  in the form

$$C_x^\pi(\ell) = \left\{ \theta : |\theta - m(x)| \leq \ell \sigma_\tau \right\}$$

for some  $\ell$  which does not depend on  $x$ . In fact, the Bayes rule chooses  $\ell$  to minimize the posterior expected loss

$$L(C_x^\pi(\ell) | x) = \bar{S}(\ell \sigma_\tau) + P(|Z| > \ell),$$

where  $Z \sim N_p(0, 1)$ .

If we now apply Theorem 3.2, we get a sufficient condition: the Bayes sets are nontrivial if there exists  $\ell_0$  such that either

$$(3.7) \quad \sigma_\tau \bar{S}'(\ell \sigma_\tau) - f_Z(\ell) < 0 \quad \text{for } 0 < \ell < \ell_0$$

or

$$(3.8) \quad \sigma_\tau \bar{S}'(\ell \sigma_\tau) - f_Z(\ell) > 0 \quad \text{for } \ell > \ell_0,$$

where  $f_Z(\cdot)$  is the density of  $|Z|$  and  $f_Z(t) = t^{p-1} e^{-t^2/2} / [2^{(p/2)-1} \Gamma(p/2)]$ .

We can now apply these formulas to the size functions in (3.6) to obtain a number of interesting results.

**Theorem 3.3.** *Let  $X | \theta \sim N_p(\theta, \sigma^2 I)$ ,  $\theta \sim N_p(\mu, \tau^2 I)$*

- a. *For the rational size function (3.6i), the Bayes sets are nontrivial for any  $a > 0$ .*
- b. *For the probability size function (3.6ii), the Bayes sets are nontrivial for all  $Z$  if and only if  $a < \sigma^{-1}$ .*

For the exponential size function (3.6iii), the Bayes sets are nontrivial if  $p = 1$ . For  $p \geq 2$ , the Bayes sets are nontrivial for all  $Z$  if and only if  $a < \sigma^{-1}$ .

**Proof** a. Condition (3.8) is obviously satisfied, as the rational size function has a heavier tail than  $f_Z(\cdot)$ , which decreases exponentially fast.

b. For the probability size function we have

$$L\left(C_X^\pi(\ell) \mid x\right) = P\left(|Z| \leq a\ell\sigma_\tau\right) + P\left(|Z| > \ell\right),$$

which is clearly negative if  $a\sigma_\tau < 1$ . Thus, since  $\sigma_{\tau-1} > \sigma-2$ , part b is proved.

c. From conditions (3.7)–(3.8), we will have nontrivial Bayes sets if

$$\frac{\partial}{\partial \ell} L\left(C_X^\pi(\ell) \mid x\right) = \sigma_\tau^2 a \ell e^{-a(\ell\sigma_\tau)^2/2} - \ell^{p-1} e^{-\ell^2/2} / \left[2^{(p/2)-1} \Gamma(p/2)\right]$$

is positive for all sufficiently large  $\ell$  or negative for  $\ell$  sufficiently close to zero. If  $p = 1$ , it is straightforward to check that  $(\partial/\partial \ell)L\left(C_X^\pi(\ell) \mid x\right) < 0$  for  $\ell$  near 0, so the Bayes sets are nontrivial. If  $p \geq 2$  and  $a\sigma_\tau \geq 1$ , then it is also straightforward to check that  $L\left(C_X^\pi(\ell) \mid x\right) \geq 0$  for all  $\ell$ , so the Bayes sets are trivial by Theorem 4.1. If  $a\sigma_\tau < 1$ , then it can be shown that (3.9) is positive for all sufficiently large  $\ell$  and hence, the Bayes sets are nontrivial. The condition follows since  $\sigma_{\tau-1} > \sigma-2$ .  $\square$

The three size functions of (3.6) have been chosen because they span a range of possibilities. The relationship between the size function and the density function is of utmost importance in determining the behavior of the Bayes sets. With the size functions in (3.6), we have tails that can be heavier than normal (rational size), lighter than normal (probability size) and equal to normal (exponential size). It follows from the previous developments that the size function must have a tail heavier than the posterior density in order for the Bayes sets to be nontrivial. In less technical words, if the penalty for large size increases rapidly enough (to its upper bound), or increases slowly enough (from its lower bound), the Bayes sets will be trivial.

### 3.3. Bounds on Bayes Sets

In this section, we still concentrate on the normal case using the size functions defined in (3.6). For these cases, we find that as the value of  $a$  changes, there is a smallest, nonempty Bayes set.

However, there is no largest bounded Bayes set (that is, the maximum volume of a Bayes set is infinity.)

To be explicit, for a specified value of  $a$  we write the Bayes set

$$C_x^\pi[\ell(a)] = \left\{ \theta : |\theta - m(x)| \leq \ell(a)\sigma_\tau \right\}.$$

For each loss given in (3.6), there is a unique value of  $\ell(a)$  that describes the Bayes set. We now find the range of this function as  $a$  varies. Consider first the rational size function.

**Theorem 3.6.** *If  $X|\theta \sim N_p(\theta, \sigma^2I)$  and  $\theta \sim N_p(\mu, \tau^2I)$  the Bayes set under rational size is such that  $\ell(a)$  satisfies*

- i)  $\max_a \ell(a) = \infty$
- ii)  $\min_a \ell(a) = \ell^* > 0$ .

**Proof.** We can write the posterior loss of  $C_x^\pi[\ell(a)]$  as

$$L\left(C_x^\pi[\ell(a)] \mid x\right) = \frac{(\ell(a)\sigma_\tau)^p}{a + (\ell(a)\sigma_\tau)^p} + P\left(\chi_p^2 > \ell^2(a)\right).$$

As  $a \rightarrow 0$ ,  $L\left(C_x^\pi[\ell(a)] \mid x\right)$  can remain negative (as it must, by Theorem 3.3) only if  $\ell \rightarrow \infty$ . This establishes that  $(\partial/\partial a)L\left(C_x^\pi[\ell(a)] \mid x\right) < 0$ . This implies that the minimum value of  $\ell(a)$  is attained for  $0 < a < \infty$ , establishing part ii). □

It is straightforward to verify that the minimum value of  $\ell(a)$  is independent of  $\sigma_\tau$ . Although analytic determination of this minimum is difficult, a numerical algorithm can be used. The following

table gives some selected values of  $\ell^* = \min_a \ell(a)$ :

Table 3.1. For dimension  $p = 1, 2, \dots, 7$ , values of the minimum length  $\ell^*$  and corresponding coverage of conjugate Bayes sets against the rational loss.

$p$	1	2	3	4	5	6	7
$\ell^*$	1.90	2.07	2.21	2.29	1.40	1.40	1.29
Posterior coverage	.94	.88	.82	.74	.31	.08	.02

Of course, the posterior coverage of the Bayes set  $C_X^\pi(\ell^*)$  is the minimum probability, and corresponds to the usual frequentist coverage when the prior is taken to be Lebesgue measure. Table 3.1 shows that for small values of  $p$ , the rational size function dictates a rather large coverage probability, and hence a restricted range of Bayes sets. As  $p$  increases, however, the range of the coverage probabilities increases.

We next turn to the probability size function.

**Theorem 3.5.** For  $X|\theta \sim N_p(\theta, \sigma^2I)$  and  $\theta \sim N_p(\mu, \tau^2I)$ , under probability size function, with  $a < \sigma^{-1}$ , the Bayes set is associated with

$$\ell(a) = \left[ \frac{-2p \log(a\sigma_\tau)^2}{1 - (a\sigma_\tau)^2} \right]^{\frac{1}{2}}$$

and  $\min_{a < \sigma_\tau^{-1}} \ell(a) = p^{1/2}$  and  $\max_{a < \sigma_\tau^{-1}} \ell(a) = \infty$ .

Thus the probability size function also results in an unbounded range of Bayes sets. As before, the coverage probability of the smallest Bayes set,  $C_X^\pi(p^{1/2})$ , is independent of  $\sigma_\tau$ . However, for this size function the range of the Bayes sets is wider for small  $p$ , as evidenced in Table 3.2. Moreover, it can be established analytically that for this loss

$$\min_a P\left(\theta \in C_X^\pi[\ell(a)] \mid x\right) \geq \frac{1}{2} \quad \text{for any } p.$$

Table 3.2. For selected dimensions  $p$ , values of the coverage of the minimum length Bayes set against the probability size loss.

$p$	1	3	5	10	20	40	100
Posterior coverage	.69	.61	.59	.56	.54	.50	.50

For the exponential size function, the results are similar to those for the probability size function, and will only be summarized.

**Theorem 3.6.** For  $X|\theta \sim N_p(\theta, \sigma^2 I)$ ,  $\theta \sim N_p(\mu, \tau^2 I)$  and the expo size function, when  $a < \sigma^{-1}$  and  $p > 2$ , the Bayes set is associated with

$$(p-2)^{1/2} \leq \ell(a) < \infty .$$

Again, as with the probability size function, it can be established that, as lone as  $p > 2$  and  $a < \sigma_\tau^{-1}$ , the minimum coverage probability of the Bayes set is at least 1/2. In fact, an analytic bound on the length is

$$\min_a \ell(a) \geq \min_a \sqrt{\frac{p-2}{1-a^2}} = \sqrt{p-2} .$$

The coverage probabilities of the Bayes sets of minimum volume are given in Table 3.3. They indicate that the minimum length is presumably much larger than  $\sqrt{p-2}$ .

Table 3.3. For selected dimensions  $p$ , values of the minimum length  $\ell^*$  and coverage of the Bayes set (4.15) against the loss (4.18).

$p$	3	4	5	6	7	8	9
$\ell^*$	2.24	2.66	3.00	3.27	3.51	3.74	3.96
Posterior coverage	.83	.87	.89	.90	.91	.92	.93
$p$	10	11	12	13	14	15	20
$\ell^*$	4.15	4.32	4.49	4.65	4.80	4.96	5.65
Posterior coverage	.93	.93	.93	.94	.94	.94	.96

#### 4. Relations Between the Loss and the Distribution

When we consider a given loss,

$$(4.1) \quad L(\theta, C) = S \text{vol}(C) + \mathbb{1}(\theta \notin C) ,$$

one important problem is to determine how much this loss depends on the distribution. For instance, we can look at the conditions for admissibility or minimaxity of different sets, and see how these conditions depend on the distribution. Another subject of interest is the value of the minimal coverage probability over a class of losses, as treated in the previous section. We can ask if this lower bound depends on the distribution.

In general, there will be a relationship between the loss and the distribution. We can then try to restrict ourselves to the consideration of “robust” or “universal” losses. Although this notion is not easy to define, it is certain that we cannot define a loss that will be valid for any kind of problem. This problem is not unique to set estimation, but also is the case for point estimation problems. In point estimation problems, it is well known that different combinations of loss functions can give different “optimal” answers for the same distribution.

We show in Section 4.1 that a natural loss exists for bounded parameter spaces. This loss can be called “universal” in the sense that all Bayes sets are nontrivial, independent of the distribution of the observations. In the more general case of unbounded parameter spaces, some of the losses we have considered will not behave properly for more general models. For instance, the size function  $S^*(C) = P(T_\nu \in C)$  gives  $\emptyset$  as a minimax set if the observations have a t distribution with  $n < \nu$  degrees of freedom. The rational loss thus appears as the most appealing candidate.

##### 4.1. Bounded Parameter Space

We now consider the case where  $X$  has a distribution depending on a parameter  $\theta \in \Theta$  with  $\text{vol}(\Theta) = \Omega < +\infty$  (the volume being determined by Lebesgue measure). Example 2.3 is a particular case of this situation, and we still see that some of the results in Example 2.3 for the binomial distribution can be generalized to distributions with bounded parameter spaces. In fact, a natural loss in this setup is

$$(4.2) \quad L(\theta, c) = \frac{\text{vol}(C)}{\Omega} + \mathbb{1}(\theta \notin C),$$

as the volume function is increasing, bounded by one and satisfies  $L(\theta, \emptyset) = L(\theta, \Theta) = 0$ . Furthermore, the following theorem shows that the loss is independent of the distribution of the observations.

**Theorem 4.1.** *For any distribution  $f(x|\theta)$  and any prior distribution  $\pi$  on  $\Theta$ , the Bayes set is*

$$C_x^\pi = \left\{ \theta : \pi(\theta|x) \geq \frac{1}{\Omega} \right\}.$$

**Proof.** In the case of the loss function (4.2), we have  $S'(t) = 1/\Omega$ . Thus, for a set of the form  $C_x^\pi(k) = \{\theta : \pi(\theta|x) \geq k\}$  it follows from Lemma 2.1 that

$$(4.3) \quad \frac{\partial}{\partial k} \left( L(\theta, C_x^\pi(k)) \right) = \left[ k - \frac{1}{\Omega} \right] \int_{\{\pi(\theta|x) = k\}} \frac{ds}{|\nabla \pi(\theta|x)|}.$$

The integral in (4.3) is always nonnegative. Furthermore, this integral has to be positive for some values of  $k$ , otherwise  $\pi$  would not be a density on  $\Theta$ . Therefore, as  $k$  increases from 0,  $\frac{\partial}{\partial k} \left( L(\theta, C_x^\pi(k)) \right)$  changes sign only once from negative to positive. This implies that  $k = 1/\Omega$  gives a minimum of  $L(\theta, C_x^\pi(k)|x)$ , for any distribution  $f(x|\theta)$ .  $\square$

**Theorem 4.2.** *For any prior distribution,  $\pi(\theta)$ , the Bayes set associated with the loss function (4.2),  $C_x^\pi$ , is nontrivial.*

**Proof.** We first prove that the posterior distribution cannot be uniform. If  $f(x|\theta)$  is the density of  $X$  with respect to a measure  $\nu$ , we have  $\pi(\theta|x) = f(x|\theta)\pi(\theta) / \int_{\Theta} f(x|\theta)\pi(\theta)d\theta$ . If  $\pi(\theta|x) = 1/\Omega$  for almost every  $x$  (with respect to  $\nu$ ), then  $f(x|\theta)\pi(\theta) = (1/\Omega) \int_{\Theta} f(x|\theta)\pi(\theta)d\theta$ . Integrating with respect to  $x$  gives

$$\begin{aligned} \pi(\theta) &= \pi(\theta) \int_{\mathfrak{X}} f(x|\theta)d\nu(x) = \frac{1}{\Omega} \int_{\mathfrak{X}} \int_{\Theta} f(x|\theta)\pi(\theta)d\theta d\nu(x) \\ &= \frac{1}{\Omega} \int_{\Theta} \int_{\mathfrak{X}} f(x|\theta)d\nu(x)\pi(\theta)d\theta = \frac{1}{\Omega}. \end{aligned}$$

But if  $f(x|\theta)$  actually depends on  $\theta$ , the posterior distribution associated with  $\pi(\theta) = 1/\Omega$  is not uniform. Thus we have a contradiction.

As  $\int_{\Theta} \pi(\theta|x)d\theta = 1$  and  $\pi(\theta|x) \neq 1/\Omega$ , it follows that the sets  $\{\theta : \pi(\theta|x) > (1/\Omega)\}$  and  $\{\theta : \pi(\theta|x) < (1/\Omega)\}$  are not equivalent to the empty set or the entire space.  $\square$

We can deduce from Theorem 4.2 that  $\phi$  is not Bayes for (4.2) and, therefore, less likely to be admissible or minimax.

#### 4.2. The Rational Loss

For the general case of unbounded spaces, we can no longer use the “natural” loss of Section 4.1. Closer examination of the bounded parameter space case helps us understand why the loss seems “universal.” First, the linearity of the loss dictates that the tail behavior of the linear loss is appropriate. In fact, since the derivative of linear loss is constant, it follows that, as  $x \rightarrow \infty$ , every density function  $f(x|\theta)$  will eventually be below this constant. This is the fact that leads to the universality of linear loss for bounded parameter spaces.

Secondly, only in bounded parameter spaces will the linear loss be coherent in the sense that  $S(0) = 0$ ,  $S(\infty) = 1$ , and  $0 < S(v) < 1$  for  $0 < v < \infty$ . Thus, to try to carry these conditions over to the unbounded parameter space problem, we must consider a loss function whose derivative is not too small. A reasonable candidate is the family of ‘rational’ losses, since they avoid trivial Bayes sets for conjugate priors and for  $p \geq 5$ , the minimum coverage probability phenomenon is not of consequence in the normal case. An additional sufficient condition for nontriviality is given by the following result.

**Theorem 4.3.** *If  $\{\theta: \pi(\theta|x) > (1/a)\}$  has positive Lebesgue measure, the Bayes rule  $C_x^\pi = \{\theta: \pi(\theta|x) > k\}$  against the rational loss function is nontrivial, where  $k = k(x) < 1/a$ .*

**Proof.** See Casella, Hwang and Robert (1992).

That  $\{\pi(\theta|x) > (1/a)\} = \phi$  does not necessarily imply that the Bayes rule is  $\phi$ , as illustrated by the following example.

**Example 4.1.** Suppose  $X \sim N_p(\theta, I)$  and  $\theta \sim N_p(0, \tau^2 I)$ . Then  $\pi(\theta|x)$  is  $N_p(\eta x, \eta I)$  and the Bayes set is  $C_x^\pi = \{\theta: |\theta - \eta x|^2 \leq c\}$  where  $\eta = \tau^2/(\tau^2+1)$ . If  $a$  is such that  $(1/a) > (2\pi\eta)^{-p/2}$ , it follows that  $\{\pi(\theta|x) > (1/a)\}$  is empty. However, as seen in Section 3.3, there is still a nonempty Bayes set.

## 5. Conclusion

This paper shows how difficult it is to treat the set estimation problem in a decision-theoretic way, since it points out that many losses, even the nonlinear losses like in (1.2), suffer from defects reproached to the original linear loss. It also provides some guidelines in the selection of the size function  $S$ , in particular through the “coherency conditions”  $S(0) = 0$ ,  $S(\infty) = s$ . Unfortunately, it appears that the distribution of the observations,  $f(x|\theta)$ , also plays an important role in the apparition of paradoxes and undesirable features, like the existence of trivial Bayes sets or the minimaxity of  $\phi$ . This implies further robustness studies for specific families of losses and distributions.

Despite these limitations, we deeply recommend a decision-theoretic approach to confidence set estimation and this for three reasons. First, an independent approach separates the testing imperatives from the confidence imperative and allows us to evaluate the confidence procedures for themselves. Second, in a related way, it compels the decision-maker to evaluate the consequences of his/her actions related to the confidence procedure, thus repositions confidence procedures at the core of the inferential process instead of presenting them as secondary tools. Last, by allowing the use of classical decision-theoretic notions like admissibility, minimaxity, etc., this approach brings the decision-maker on a sounder ground since he/she can compare estimators and look for an optimal solution, as is the case in point estimation, instead of borrowing from resting perspectives or devising *ad hoc* optimality notions.

A deeper purpose of this paper is also to call for a reflection on the actual goals of confidence set estimation. Indeed, although the previous sections show that a decision-theoretic approach is possible for a careful choice of the loss function (1.3), the true purpose of decision theory is still to provide the decision-maker with a tool to help him to make rational/coherent decisions. Therefore, following De Groot (1970), Berger (1985) or Lindley (1965), this implies that the loss function he/she uses in this process is obtained from a corresponding utility function. It is yet unclear how this utility function can be constructed in practice, given the available choice of acceptable loss functions, but it is clear that the loss function should not be chosen for its mathematical properties or even for its compliance

with an “intuitive” behavior.

Borrowing from Berger (1985, Chap. 2), we propose the following approximation for the determination of the size function, assuming it belongs to a parametrized family  $S_a$ , like the rational size  $S_a(v) = v/a + \sigma$ . If the decision-maker is in fact able to compare the consequences of a large set with those of a smaller set, he can produce a sequence of volumes  $v_1, \dots, v_m$  and a corresponding sequence of weights  $\rho_2, \dots, \rho_m$  such that  $S_a(v_2) = \rho_2 S_a(v_1), \dots, S_a(v_m) = \rho_m S_a(v_2)$ . An approximate value of  $a$  can then be estimated based on this information.

### References

- Berger, J. O. (1985). *Statistical Decision Theory and Bayesian Analysis*, 2nd edition. Springer-Verlag.
- Blyth, C. R. and Hutchinson, D. W. (1961). Tables of Neyman-Shortest Confidence Intervals for the Binomial Parameter. *Biometrika* **47**, 381-391.
- Casella, G., Hwang, J. T-G., and Robert, C. P. (1992). A paradox in decision-theoretic interval estimation. *Statist. Sinica*, to appear.
- Cohen, A. and Strawderman, W. E. (1973). Admissibility implications for different criteria in confidence estimation. *Ann. Statist.* **1**, 363-366.
- De Groot, M. (1970). *Optimal Statistical Decisions*. John Wiley & Sons, New York.
- Hwang, J. T. and Casella, G. (1982). Minimax confidence sets for the mean of a multivariate normal distribution. *Ann. Statist.* **10**, 868-881.
- Hwang, J. T. and Casella, G. (1984). Improved set estimators for a multivariate normal mean. *Statistics and Decisions, Supplement Issue 1*, 3-16.
- Joshi, V. M. (1969). Admissibility of the usual confidence set for the mean of a univariate or bivariate normal population. *Ann. Math. Statist.* **40**, 1042-1067.
- Lindley, D. (1985). *Making Decisions*. John Wiley & Sons, New York.
- Meeden, G. and Vardeman, S. (1985). Bayes and admissible set estimation. *J. Amer. Statist. Assoc.* **80**, 465-471.

**Appendix**

**Theorem (Generalization of Theorem 3.1).** For  $X \sim f(x|\theta)$  and (possibly improper) prior  $\pi(\theta)$  the Bayes estimator,  $C_x^\pi$ , against the loss

$$L_{\mathbf{S}}(\theta, C) = \mathbf{S}[\text{vol}(C)] - \mathbb{I}(\theta \in C) ,$$

satisfies

$$(A.1) \quad C_x^\pi \subset \{\theta : \pi(\theta|x) \geq k\} \quad \text{or} \quad \{\theta : \pi(\theta|x) \geq k\} \subset C_x^\pi$$

for every  $k > 0$ . Furthermore, if  $C_x^\pi$  satisfies

$$\{\theta : \pi(\theta|x) \geq k_1\} \subset C_x^\pi \subset \{\theta : \pi(\theta|x) \geq k_2\}$$

where  $k_1 > k_2$  and  $\{\theta : \pi(\theta|x) = k\} = \emptyset$  for  $k_2 < k < k_1$ , then there exists  $k^* \in [k_2, k_1]$  such that

$$C_x^\pi = \{\theta : \pi(\theta|x) \geq k^*\}$$

and  $k^*$  minimizes

$$\mathbf{S}(\text{vol}\{\theta : \pi(\theta|x) \geq k\}) - \int_{\{\theta : \pi(\theta|x) \geq k\}} \pi(\theta|x) d\theta .$$

**Proof.** Since  $C_x^\pi$  is a Bayes set it minimizes the posterior loss

$$L(C_x^\pi|x) = \mathbf{S}[\text{vol}(C)] - \int_C \pi(\theta|x) d\theta .$$

If (A.1) is not satisfied, there exists  $k \geq 0$  such that

$$C_x^\pi \cap \{\theta : \pi(\theta|x) < k\} \neq \emptyset \quad \text{and} \quad (C_x^\pi)^c \cap \{\theta : \pi(\theta|x) \geq k\} \neq \emptyset ,$$

the intersections being different from zero (recall that we are working with sets defined only up to sets of Lebesgue measure zero). Therefore, there exists sets  $A_x$  and  $B_x$  such that

$$A_x \subset C_x^\pi \cap \{\theta : \pi(\theta|x) < k\}, \quad B_x \subset (C_x^\pi)^c \cap \{\theta : \pi(\theta|x) \geq k\}$$

and  $\text{vol}(A_x) = \text{vol}(B_x) > 0$ . If we now define  $C_x^* = (C_x^\pi - A_x) \cup B_x$ , it follows that

$$L(C_x^\pi|x) > L(C_x^*|x) ,$$

as  $\text{vol}(C_x^\pi) = \text{vol}(C_x^*)$  and  $\int_{A_x} \pi(\theta|x) d\theta < \int_{B_x} \pi(\theta|x) d\theta$ . Thus we have a contradiction, so  $C_x^\pi$  must satisfy (A.1). The second part of the theorem follows immediately by continuity.  $\square$