

September 1968

SOME CONDITIONAL ESTIMATION PROBLEMS  
WITH APPLICATIONS TO ESTIMATING  
THE PROBABILITY OF MISCLASSIFICATION\*

by

Marilyn Sorum and Robert J. Buehler

Technical Report No. 112

University of Minnesota  
Minneapolis, Minnesota

\*This research was supported by National Science Foundation Grant GP-6859.

## 1. INTRODUCTION

Let  $x$  denote observations arising from a distribution  $f(x; \theta)$  where  $\theta$  is unknown, and let  $R(x)$  denote some known function of the data  $x$ . We will be concerned with the estimation from  $x$  of functions of the form  $\phi(R(x), \theta)$ . The problem arises in the theory of classification where  $\phi$  denotes the probability that a classification rule  $R$  based on data  $x$  will misclassify a future observation from a specified population characterized by an unknown parameter  $\theta$ . A natural way to construct an estimator of  $\phi$  would be to put  $\hat{\phi}(x) = \phi(R(x), \hat{\theta}(x))$  where  $\hat{\theta}(x)$  is, say, the maximum likelihood estimator of  $\theta$ . We shall however approach the problem after the fashion of the standard Cramér-Rao and Rao-Blackwell theory, considering the conditional variance of estimators which are conditionally unbiased, given  $R(x)$ . Two examples are given, one involving a classification problem for normal distributions.

## 2. THE CRAMÉR-RAO BOUND

The usual Cramér-Rao theory extends in a straightforward way to include the present case. To simplify the presentation, we will treat the case of  $n$  observations from a continuous distribution, putting  $\underline{x} = (x_1, \dots, x_n)$ . We further suppose that  $R(\underline{x})$  is such that we can find  $(n-1)$  coordinates  $y_1, \dots, y_{n-1}$  giving a "smooth" transformation (one-to-one with Jacobian existing and nowhere equal to 0 or  $\infty$ ) from  $\underline{x}$  to  $\underline{y}$ , where  $\underline{y} = (y_1, \dots, y_n)$ ,  $y_n = R(\underline{x})$ .

Theorem 1. Let  $\phi(R, \theta)$  be a function such that  $\partial\phi/\partial\theta$  exists. Let  $g_\theta(y_1, \dots, y_{n-1}|y_n)$  denote the conditional density of  $y_1, \dots, y_{n-1}$  given  $y_n$ , and assume this density satisfies the usual regularity conditions (see, for example, Lehmann [6]). Let  $\hat{\phi}(\underline{x})$  be any conditionally unbiased estimator (i.e.,  $E(\hat{\phi}(\underline{x})|Y_n = y_n) = \phi(y_n, \theta)$ ). Then  $\hat{\phi}$  satisfies

$$(1) \quad \text{Var}(\hat{\phi}(\underline{x})|Y_n = y_n) \geq \{E(S^2|Y_n = y_n)\}^{-1}$$

where

$$(2) \quad S = \partial \log g_\theta / \partial \theta = (\partial \log g_\theta / \partial \theta)(\partial\phi/\partial\theta)^{-1}.$$

We suppress the proof, which follows the usual lines. Superficially it appears that the bound (1) depends on the arbitrary choice of coordinates  $y_1, \dots, y_{n-1}$ . But this is not actually so, for let  $z_1, \dots, z_{n-1}$  be any other choice and let  $g_\theta^*$  and  $S^*$  be the analogues of  $g_\theta$  and  $S$ . Then  $g_\theta^* = Jg_\theta$ , where  $J = |\partial(y_1, \dots, y_{n-1})/\partial(z_1, \dots, z_{n-1})|$ , and  $S^* = S$ , so that  $S^*$  and  $S$  have the same second moment.

It is possible to avoid the regularity conditions on  $g_\theta$  and on  $\phi(R, \theta)$  by the method of Chapman and Robbins [3]. We will state the result without proof.

Theorem 2. In the notation of Theorem 1, define

$$(3) \quad A = A(\phi, \theta, h, r) = h^{-1}[\phi(r, \theta+h) - \phi(r, \theta)]$$

$$(4) \quad B = B(\phi, \theta, h, r) = A^{-2} h^{-2} \left\{ \left[ \frac{g_{\theta+h}(y_1, \dots, y_{n-1}|r)}{g_\theta(y_1, \dots, y_{n-1}|r)} \right]^2 - 1 \right\}.$$

Then for any conditionally unbiased estimator  $\hat{\phi}(\underline{x})$  of  $\phi(R, \theta)$ ,

$$(5) \quad \text{Var}(\hat{\phi}|r) \geq \{\inf E[B(\phi, \theta, h, r)|r]\}^{-1}$$

where the infimum is taken over all  $h \neq 0$  such that both  $\theta$  and  $\theta+h$  lie in the range of possible  $\theta$  values.

### 3. RAO-BLACKWELL THEORY

In order to give the desired modification of the Rao-Blackwell theory, we will require definitions of conditional sufficiency and conditional complete sufficiency.

Definition 1. A statistic  $T(x)$  is called conditionally sufficient for  $\theta$  given  $R(x)$  if  $(R(x), T(x))$  is sufficient for  $\theta$ .

To discuss completeness, we will consider separately the discrete case and the absolutely continuous case.

Definition 2 (discrete case). A statistic  $T(x)$  is called a conditionally complete sufficient statistic for  $\theta$  given  $R(x)$  if for each value  $r$  of  $R$

$$(6) \quad \sum_t f(t) P_{\theta}(T=t|R=r) = 0 \quad \text{for all } \theta$$

implies

$$(7) \quad f(t) = 0 \quad \text{for all } t \text{ such that } P_{\theta}(T=t|R=r) > 0.$$

Definition 2 (absolutely continuous case). Let  $h_{\theta}(\cdot)$  be the density of a statistic  $R(x)$  and let  $g_{\theta}(\cdot|r)$  be the conditional density of a statistic  $T(x)$  given  $R(x) = r$ , defined for all  $r$  such that  $h_{\theta}(r) > 0$ .  $T(x)$  is called conditionally complete sufficient given  $R$  if for each  $r$  such that  $h_{\theta}(r) > 0$ ,

$$(8) \quad \int f(t) g_{\theta}(t|r) dt = 0 \quad \text{for all } \theta$$

implies

$$(9) \quad f(t) = 0 \quad \text{a.e. } P_{\theta}^{T|R=r}.$$

It is more or less evident from Definition 1 that every sufficient statistic is conditionally sufficient (see Bahadur [2], Theorem 6.4, for a rigorous treatment). It is also true that completeness implies conditional completeness. We give a proof for the discrete case only.

Theorem 3. Let  $P_\theta$  be a family of probability measures defined on a discrete space  $\{x\}$  assigning nonzero probability to each point for each value of  $\theta$ . If  $T(x)$  is a complete sufficient statistic for  $\theta$ , then  $T(x)$  is a conditionally complete sufficient statistic for  $\theta$  given any statistic  $R(x)$ .

Proof. Let  $\mathcal{J}$ ,  $\mathcal{R}$  and  $\mathcal{S}$  denote the range of  $T(x)$ ,  $R(x)$  and  $(T(x), R(x))$  respectively. Every point of  $\mathcal{J}$  and  $\mathcal{R}$  has nonzero probability for every  $\theta$ . The same is true of  $\mathcal{S}$  although  $\mathcal{S}$  is not necessarily the direct product of  $\mathcal{J}$  and  $\mathcal{R}$ . For any point  $(t, r)$  in  $\mathcal{S}$  both  $P_\theta(T=t|R=r)$  and  $P_\theta(R=r|T=t)$  are defined and nonzero for every  $\theta$ . By sufficiency of  $T$ , the latter is the same for all  $\theta$  and therefore can be denoted by  $C(r, t)$ . Now consider any fixed  $r$  in  $\mathcal{R}$ . We can write

$$(10) \quad \sum_t f(t)P_\theta(T=t|R=r) = \frac{1}{P_\theta(R=r)} \sum_t f(t)C(r, t)P_\theta(T=t).$$

If we assume (6) holds, then the above expression equals zero for all  $\theta$ . The assumed completeness of  $T(x)$  then implies  $f(t)C(r, t) = 0$  for all  $t$  in  $\mathcal{J}$ . Since  $C(r, t) \neq 0$  for  $(r, t)$  in  $\mathcal{S}$ , we have  $f(t) = 0$  whenever  $(r, t)$  in  $\mathcal{S}$ , that is, whenever  $P_\theta(T=t|R=r) > 0$ , so that (7) holds.

We now state without proof the modification of the standard Rao-Blackwell Theory.

Theorem 4. Let  $\hat{\phi}$  be any estimator of  $\phi$  such that  $E_{\theta}(\hat{\phi}|R = r) = \phi(r, \theta)$ . Let  $T$  be conditionally sufficient for  $\theta$  given  $R$ . Define  $\tilde{\phi} = E(\hat{\phi}|R, T)$ . Then  $\tilde{\phi}$  is also a conditionally unbiased estimator of  $\phi$  given  $R$ , and  $\text{Var}(\tilde{\phi}|R) \leq \text{Var}(\hat{\phi}|R)$ .

Lemma 1. Let  $\hat{\phi}(x)$  and  $\psi(x)$  be two conditionally unbiased estimators of  $\phi(R, \theta)$  based on a conditionally complete sufficient statistic. Then  $\hat{\phi}(x) = \psi(x)$  a.e.  $P_{\theta}^X|R$ .

Theorem 5. Let  $T$  be conditionally complete sufficient for  $\theta$ . Let  $\phi(R, \theta)$  be any quantity for which a conditionally unbiased estimator given  $R$  exists. Then  $\phi(R, \theta)$  has a unique (a.e.) conditionally UMVU estimator which is a function of  $T$  and  $R$ .

Thus when a conditionally unbiased estimator is known, we can find a conditionally UMVU estimator by calculating the conditional expectation given both  $R$  and  $T$ .

#### 4. EXAMPLES

4.1 Example 1. Let  $x_1, \dots, x_n$  be a sample from  $N(\mu, 1)$  and let  $R(x) = \sum c_i x_i$ . To avoid a degenerate case we assume not all  $c_i$  are equal. If  $y_j = x_j$ ,  $j = 1, \dots, n-1$ , and  $y_n = R(x)$ , then the conditional distribution of  $y_1, \dots, y_{n-1}$  given  $y_n$  is multinormal with mean

$$(11) \quad \bar{y} = d^{-1}rc + \mu(1-d^{-1}bc)$$

and covariance matrix

$$(12) \quad C = I_{n-1} - d^{-1}cc'$$

where

$$(13) \quad b = \sum_1^n c_i, \quad d = \sum_1^n c_i^2, \quad \underline{c} = (c_1, \dots, c_{n-1})', \quad \underline{1} = (1, \dots, 1)'$$

We find

$$(14) \quad \partial \log g_\mu(y_1, \dots, y_{n-1} | r) / \partial \mu = (\underline{1} - d^{-1} b \underline{c})' C^{-1} (\underline{y} - \underline{v})$$

and

$$(15) \quad E(S^2 | r) = V \left( \frac{\partial \varphi}{\partial \mu} \right)^{-2}, \quad V = (\underline{1} - d^{-1} b \underline{c})' C^{-1} (\underline{1} - d^{-1} b \underline{c}).$$

If we specialize to  $n = 3$ ,  $R(x) = x_1 + x_2 + 2x_3$ , then any conditionally unbiased estimator of  $\varphi(R, \mu)$  satisfies

$$\text{Var}(\hat{\varphi} | r) \geq 3(\partial \varphi / \partial \mu)^2.$$

This bound may or may not be achievable depending on the form of  $\varphi(R, \mu)$ . It is, for example, achieved when  $\varphi(R, \mu) = R\mu$  if we take  $\hat{\varphi} = R\hat{\mu}$  with  $\hat{\mu} = \frac{3}{2}(x_1 + x_2) - \frac{1}{2}r = x_1 + x_2 - x_3$ . The choice  $\hat{\mu} = \bar{x}$  would not give a conditionally unbiased estimator.

#### 4.2 Example 2.

4.2.1 A problem in classification theory. The present paper was in fact motivated by some problems in classification theory which are more fully discussed in [7]. Let  $S_1$  and  $S_2$  denote samples of size  $N_1$  and  $N_2$  known to come from populations  $\pi_1$  and  $\pi_2$  respectively. The samples determine a rule for classifying a future observation as either belonging to  $\pi_1$  or to  $\pi_2$ . We wish to estimate the probabilities of misclassifying an observation from  $\pi_1$  in  $\pi_2$  or vice versa. These two probabilities depend on the samples and on unknown population parameters.

Consider next the p-variate normal case with unknown mean vectors  $\underline{\mu}_1, \underline{\mu}_2$  and known and equal covariance matrices  $\Sigma$ . If

$\bar{x}_1, \bar{x}_2$  are the sample mean vectors and  $\underline{z}$  is a future observation, then the usual symmetrical classification rule is (Anderson [1])

(16) classify  $\underline{z}$  as  $\pi_1$  if  $(\bar{x}_1 - \bar{x}_2)' \Sigma^{-1} \underline{z} \geq \frac{1}{2}(\bar{x}_1 - \bar{x}_2)' \Sigma^{-1} (\bar{x}_1 + \bar{x}_2)$  and otherwise classify  $\underline{z}$  as  $\pi_2$ . For definiteness consider the probability  $P_2$  of misclassifying an observation from  $\pi_2$  into  $\pi_1$  for given values of  $\bar{x}_1, \bar{x}_2$ , which is

$$(17) P_2 = 1 - F(C)$$

where

$$(17a) C = \frac{1}{2}D + D^{-1}(\bar{x}_1 - \bar{x}_2)' \Sigma^{-1} (\bar{x}_2 - \mu_2),$$

where  $D$  is the positive square root of

$$D^2 = (\bar{x}_1 - \bar{x}_2)' \Sigma^{-1} (\bar{x}_1 - \bar{x}_2),$$

and where  $F$  is the standard normal C.D.F.  $P_2$  is a function of  $\bar{x}_1, \bar{x}_2$  and  $\mu_2$  which we wish to estimate. In this case if we attempt to estimate conditionally on  $(\bar{x}_1, \bar{x}_2)$ , then there is no conditionally unbiased estimator of  $P_2$  because  $P_2$  depends on  $\mu_2$  but the conditional distribution of  $(S_1, S_2)$  given  $(\bar{x}_1, \bar{x}_2)$  is the same for all  $(\mu_1, \mu_2)$ , owing to the sufficiency of  $(\bar{x}_1, \bar{x}_2)$ .

The theory of conditional unbiased estimation will however apply in the case where the classification rule is tested by using it on additional observations of known origin. Although this does not seem to be a widely used method, we find it mentioned, for example by Hills [4] and Lachenbruch and Mickey [5]. Let  $t_1, \dots, t_m$  be additional observations from  $\pi_2$  and let  $q$  be the proportion of these  $m$  observations which are misclassified.



Then  $q$  is clearly a conditionally unbiased estimator of  $P_2$  since its conditional distribution given  $\bar{x}_1, \bar{x}_2$  is binomial with mean equal to  $P_2$ . (The argument does not involve the normality assumption, so that the method is clearly quite general).

4.2.2 The Cramér-Rao bound. To illustrate Theorem 1 we will apply it to the problem just described. Since Theorem 1 applies only for scalar parameters we specialize Section 4.2.1 to the univariate normal case where  $\pi_1$  and  $\pi_2$  are  $N(\mu_1, 1)$  and  $N(\mu_2, 1)$ . The function  $\varphi$  of Section 2 corresponds to  $P_2$  given by (17), and in the univariate case (17) reduces to

$$(18) \quad P_2 = \begin{cases} 1-F(c) & \text{if } \bar{x}_1 > \bar{x}_2 \\ F(c) & \text{if } \bar{x}_1 < \bar{x}_2 \end{cases}$$

where

$$(19) \quad c = \frac{1}{2} \bar{x}_1 + \frac{1}{2} \bar{x}_2 - \mu_2.$$

The conditioning variate  $y_n$  of Theorem 1 corresponds to the two sample means  $\bar{x}_1, \bar{x}_2$ . For  $y_1, \dots, y_{n-1}$  we may choose  $N_1 + N_2 + m - 2$  suitable coordinates such that, together with  $\bar{x}_1, \bar{x}_2$ , these coordinates are in one-to-one correspondence with the sample  $(x_{11}, \dots, x_{1N_1}, x_{21}, \dots, x_{2N_2}, t_1, \dots, t_m)$ . Our arbitrary choice is to delete  $x_{1N_1}$  and  $x_{2N_2}$  from the above array to obtain  $y_1, \dots, y_{n-1}$ . The desired conditional distribution is found to have mean vector  $(\bar{x}_1, \dots, \bar{x}_1, \bar{x}_2, \dots, \bar{x}_2, \mu_2, \dots, \mu_2)$  where the blocks are of length  $N_1-1, N_2-1$  and  $m$  respectively. Since  $\mu_2$  appears only in the last  $m$  positions and since the corresponding variates are (unconditionally and conditionally) independent of the others we find

$$(20) \quad \partial \log g_{\mu_2} / \partial \mu_2 = \sum_1^m (t_i - \mu_2).$$

It follows that

$$(21) \quad E(S^2 | \bar{x}_1, \bar{x}_2) = \left( \frac{\partial \varphi}{\partial \mu_2} \right)^{-2} E \left\{ \sum_{i=1}^m (t_i - \mu_2)^2 | \bar{x}_1, \bar{x}_2 \right\}$$

$$= m \left( \frac{\partial \varphi}{\partial \mu_2} \right)^{-2}.$$

Whether  $\bar{x}_1$  be greater than or less than  $\bar{x}_2$  we find

$$(22) \quad \left( \frac{\partial \varphi}{\partial \mu_2} \right)^2 = f^2(c)$$

where  $f$  is the standard normal density function. Thus the Cramer-Rao bound is

$$(23) \quad \text{Var}(\hat{\varphi} | \bar{x}_1, \bar{x}_2) \geq \frac{1}{m} f^2(c).$$

4.2.3 The conditional UMVU estimator. The estimator  $\hat{\varphi}$  of Section 4.2.1 can be improved with respect to its conditional variance by the conditional Rao-Blackwell method. Let  $\bar{t} = \frac{1}{m} \sum t_i$ . Then  $\bar{t}$  is conditionally complete for  $\mu_2$  given  $\bar{x}_1, \bar{x}_2$ . The conditional expectation of any unbiased estimator of  $P_2$  will lead to the conditional UMVU estimator, and for convenience we take  $\hat{\varphi} = 1$  or 0 according as  $t_1$  is misclassified as not by the rule (16). The UMVU estimator  $\tilde{\varphi}$  is then just the conditional probability

$$(24) \quad \tilde{\varphi} = P \left\{ (\bar{x}_1 - \bar{x}_2) \Sigma^{-1} t_1 > \frac{1}{2} (\bar{x}_1 - \bar{x}_2)' \Sigma^{-1} (\bar{x}_1 + \bar{x}_2) | \bar{x}_1, \bar{x}_2, \bar{t} \right\}.$$

To evaluate this we first find the conditional distribution of  $t_1$  given  $\bar{t}$  to be  $N_p(\bar{t}, \frac{m-1}{m} \Sigma)$ . It follows that

$$(25) \quad \tilde{\varphi} = 1 - F \left( (1-m^{-1})^{-1/2} \left[ \frac{1}{2} D + D^{-1} (\bar{x}_1 - \bar{x}_2)' \Sigma^{-1} (\bar{x}_2 - \bar{t}) \right] \right).$$

It is reasonably straightforward to get the conditional variance of  $\tilde{\varphi}$  in terms of the bivariate normal distribution function. We find

$$(26) \text{Var}(\tilde{\varphi} | \bar{x}_1, \bar{x}_2) = F(C, C; m^{-1}) - F(C, C, 0)$$

where  $C$  is defined by (17a) and where  $F(a, b; \rho) = P(u \leq a, v \leq b)$  with  $(u, v)$  joint normal with zero means, unit variances and correlation  $\rho$ .

In the univariate case, (26) holds with  $C$  replaced by  $c$  defined in (19). Table 1 compares the variance of  $\tilde{\varphi}$  with that of  $q$  and with the Cramér-Rao bound for several values of  $c$  and  $m$ .

TABLE 1

m	c	Var q	Var $\tilde{\varphi}$	C-R Bound
5	0.0	0.0500	0.0320	0.0318
	0.5	0.0427	0.0255	0.0248
	1.0	0.0267	0.0129	0.0117
	1.5	0.0125	0.0041	0.0034
	2.0	0.0044	0.0009	0.0006
10	0.0	0.0250	0.0159	0.0159
	0.5	0.0213	0.0126	0.0124
	1.0	0.0134	0.0061	0.0059
	1.5	0.0062	0.0018	0.0017
	2.0	0.0022	0.0004	0.0003
20	0.0	0.0125	0.0080	0.0080
	0.5	0.0107	0.0062	0.0062
	1.0	0.0067	0.0030	0.0029
	1.5	0.0031	0.0009	0.0008
	2.0	0.0011	0.0002	0.0001

## REFERENCES

- [1] Anderson, T. W., An Introduction to Multivariate Statistical Analysis, Wiley, New York, 1958.
- [2] Bahadur, R. R., "Sufficiency and Statistical Decision Functions," Annals of Mathematical Statistics, Vol. 25 (1954), pp. 423-462.
- [3] Chapman, D. G. and Robbins, H., "Minimum Variance Estimation Without Regularity Assumptions," Annals of Mathematical Statistics, Vol. 22 (1951), pp. 581-586.
- [4] Hills, M., "Allocation Rules and Their Error Rates," Journal of the Royal Statistical Society, Series B, Vol. 28 (1966), pp. 1-32.
- [5] Lachenbruch, P. A. and Mickey, M. R., "Estimation of Error Rates in Discriminant Analysis," Technometrics, Vol. 10 (1968), pp. 715-725.
- [6] Lehmann, E. L., Notes on the Theory of Estimation, Statistics Dept., Univ. of California, Berkeley, 1950.
- [7] Sorum, Marilyn J., "Estimating the Probability of Misclassification," Technical Report No. 110 (1968), Statistics Dept., Univ. of Minnesota, Minneapolis.