

# SMALL AREA ESTIMATION FOR THE DISTRIBUTION OF PARAMETERS

Michael P. Cohen, National Center for Education Statistics\*  
555 New Jersey Avenue NW, Washington DC 20208-5654

**Key Words:** Borrowing strength, Direct estimates, Distance, Ensemble

**Abstract:** If one wants to estimate a parameter for each of many small areas, one can generally improve the independent direct estimates by “borrowing strength” from the other small areas. Much research has been devoted to the situation in which one seeks to minimize the (possibly weighted) sums of the expected squared errors of the small area estimates. Thomas A. Louis, Malay Ghosh, and others have considered the contrasting situation in which the relationship among the small area parameters is of primary interest. For example, one might be interested in knowing the proportion of small areas where the high school dropout rate is above some level. The aim in such problems is to minimize the distance between the observed distribution of the “ensemble” (set) of small area estimates and the true distribution of the ensemble of parameters. In this paper we further explore the small area estimation problem when estimating the distribution of the parameters is the goal.

## 1. Introduction

Suppose we are investigating the values of a certain parameter (e.g. average income or an average measure of the level of literacy) for each of many small areas. If the goal is the best estimates of these parameters considered individually, then empirical and hierarchical Bayes techniques have been developed that improve upon naïve estimators. What if, though, we want to know which small areas have parameter values above a fixed cutoff  $C$  and which below? A different approach is required to treat problems of this type.

Louis (1984) was the first to study these small area estimation problems although Rubin (1981) had looked at the situation in another context. Ghosh (1992, 1994) built on the work of Louis, extending it to non-normal and multivariate situations. Our aim is to build on the work of these authors and, in particular, to investigate the use of loss functions that measure the distance between the distribution of the estimates and the distribution of the parameters.

For a general appraisal of small area estimation, Ghosh and Rao (1994) is highly recommended.

The very recent and interesting work of Shen

and Louis (1998) studies and compares the different approaches to small area estimation in a two-stage hierarchical setting.

The organization of this paper is as follows: This introduction is Section 1. Section 2 provides background information. Section 3 introduces the loss functions that will be employed. In Section 4 we study a simple normal model, and in Section 5 we extend the results to more general situations. Some concluding remarks are given in Section 6.

## 2. Background

Consider the estimation of  $m$  parameters  $\theta_1, \dots, \theta_m$  under squared error loss. Let  $\hat{\theta}_1^B, \dots, \hat{\theta}_m^B$  denote Bayes estimates of these parameters based on data  $\mathbf{X} = (X_1, \dots, X_m)$ . Let  $\theta_\bullet = \frac{1}{m} \sum_{i=1}^m \theta_i$  and  $\hat{\theta}_\bullet^B = \frac{1}{m} \sum_{i=1}^m \hat{\theta}_i^B$ . Then

$$E(\theta_\bullet | \mathbf{X}) = \hat{\theta}_\bullet^B$$

but

$$E \left[ \sum_{i=1}^m (\theta_i - \theta_\bullet)^2 | \mathbf{X} \right] > \sum_{i=1}^m (\hat{\theta}_i^B - \hat{\theta}_\bullet^B)^2.$$

This was shown by Louis (1984) under a normality assumption and, in general, by Ghosh (1992).

The point is that the Bayes estimates of the parameters (under squared error loss) have the same mean as the parameters themselves, but are on average less “spread out.” If we are trying to use the collection of Bayes estimates to study the distribution of the parameters, we will have the distorted view that the parameters are more concentrated about their mean than they really are. We have been discussing Bayes estimates, but *empirical* Bayes estimates face the same problem.

In the context of small area estimation, the  $\theta_i$  are parameters associated with small area  $i$ , say mean household income. If we use the  $\hat{\theta}_i^B$  to study the  $\theta_i$ , we will underestimate the diversity in the parameters.

Louis (1984) tackled this problem by investigating the class of estimators  $\tilde{\theta}_i$  that satisfy

$$E(\theta_\bullet | \mathbf{X}) = \tilde{\theta}_\bullet \quad \text{where } \tilde{\theta}_\bullet = \frac{1}{m} \sum_{i=1}^m \tilde{\theta}_i$$

and

$$\mathbb{E} \left[ \sum_{i=1}^m (\theta_i - \theta_*)^2 \mid \mathbf{X} \right] = \sum_{i=1}^m (\tilde{\theta}_i - \tilde{\theta}_*)^2.$$

He still used squared error loss but it was minimized subject to these constraints. The constraints force a match on the first two moments between the distribution of the estimates and the distribution of the parameters.

In giving a theoretical basis to his work, Louis (1984, Subsection 2.2) introduced the notion of a general loss function operating on the empirical distributions of the parameter estimates and the parameters. Our investigation will be based on such loss functions; they are described in the next section.

### 3. Loss Functions

Given  $m$  parameters  $\theta_1, \dots, \theta_m$ , define the function

$$G_m(t) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}(\theta_i \leq t) \quad (3.1)$$

where  $\mathbb{I}(\cdot)$  is 1 when its argument is true and 0 otherwise. We can regard  $G_m$  as the empirical distribution function of the parameters. From a Bayesian point of view, the parameters are random variables. It should be noted, however, that the parameters will generally *not* be identically distributed and maybe not independent.

Let  $\hat{G}_m$  be an estimator of  $G_m$ . For example, given  $m$  estimates  $\hat{\theta}_1, \dots, \hat{\theta}_m$  of  $\theta_1, \dots, \theta_m$  respectively, one could estimate  $G_m$  by

$$\hat{G}_m(t) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}(\hat{\theta}_i \leq t), \quad (3.2)$$

but we do not require estimators of  $G_m$  to be of the form (3.2). If we want to study the distribution of the  $\theta_i$ , we would like to find an estimate  $\hat{G}_m$  that is close, in some sense, to  $G_m$ . In other words, we would like  $\|\hat{G}_m - G_m\|$  to be small where  $\|\cdot\|$  is a distance function or metric. Examples of such distance functions include

$$\begin{aligned} & \|\hat{G}_m - G_m\|_{W,j} \\ &= \int_{-\infty}^{\infty} |\hat{G}_m(t) - G_m(t)|^j dW(t), \end{aligned} \quad (3.3)$$

$$\begin{aligned} & \|\hat{G}_m - G_m\|_{\mathbf{t},\mathbf{w},j} \\ &= \sum_{\ell=1}^L w_\ell |\hat{G}_m(t_\ell) - G_m(t_\ell)|^j, \end{aligned} \quad (3.4)$$

and

$$\begin{aligned} & \|\hat{G}_m - G_m\|_{\infty} \\ &= \max_{-\infty < t < \infty} |\hat{G}_m(t) - G_m(t)|. \end{aligned} \quad (3.5)$$

In (3.3),  $j > 0$  and  $W(t)$  is a weight function that we can choose to give more weight to ranges of parameter values in which we are especially interested. In (3.4),  $j > 0$  and the  $\mathbf{w} = (w_1, \dots, w_L)$  are weights attached to the points  $\mathbf{t} = (t_1, \dots, t_L)$ . If we adopt a general definition of integral, the second distance function is just a special case of the first. An even more special case is

$$\|\hat{G}_m - G_m\|_{t_0,j} = |\hat{G}_m(t_0) - G_m(t_0)|^j$$

that considers only a single point in the space of parameter values. For example, if  $\theta_i$  corresponds to average household income in small area  $i$  and  $t_0 = \$25,000$ , then  $|\hat{G}_m(t_0) - G_m(t_0)|$  measures how close we are in estimating the proportion of small areas with average household incomes less than or equal to \$25,000.

The distance function (3.5) is of great interest but difficult to work with analytically. There are, of course, other distance functions one might want to consider. In this paper, though, we concentrate on (3.3) with  $j = 2$ . The goal is to minimize the (conditional) expected distance given the data.

If we are presented with a distribution function estimate  $\hat{G}_m$  of the form (3.2), we can recover the set of values of the  $\hat{\theta}_i$  from the jumps in the function  $\hat{G}_m$ , but we cannot determine uniquely which small area  $i$  is associated with which jump. In fact, any one-to-one assignment of the small areas to the jumps gives rise to the same value of  $\hat{G}_m$ . Letting  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_m)$  and  $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \dots, \hat{\theta}_m)$ , Louis (1984, p. 394) suggests using a loss function of the form

$$\|\hat{G}_m - G_m\| + \epsilon \mathcal{L}(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}) \quad (3.6)$$

for some small  $\epsilon > 0$  where  $\mathcal{L}(\cdot, \cdot)$  is, for example, the sum of squared errors  $\mathcal{L}(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}) = \sum_{i=1}^m (\hat{\theta}_i - \theta_i)^2$ . The second term in the loss function is designed to force a unique assignment of the jumps in  $\hat{G}_m$  of form (3.2) to the small areas  $i$  without otherwise affecting the loss function much.

Given any estimator  $\hat{G}_m$  of  $G_m$ , not necessarily of form (3.2), we can estimate the ensemble  $\{\theta_1, \theta_2, \dots, \theta_m\}$  by

$$\begin{aligned} & \{\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_m\} \\ &= \left\{ \hat{G}_m^{-1} \left( \frac{1}{m+1} \right), \hat{G}_m^{-1} \left( \frac{2}{m+1} \right), \dots, \right. \\ & \quad \left. \hat{G}_m^{-1} \left( \frac{m-1}{m+1} \right), \hat{G}_m^{-1} \left( \frac{m}{m+1} \right) \right\}. \end{aligned}$$

We use (3.6) to determine which  $\hat{\theta}_i$  corresponds to

$$\hat{G}_m^{-1} \left( \frac{1}{m+1} \right),$$

and so forth.

In the next section, we make use of some of the loss functions described in this section to investigate a simple normal model.

#### 4. Simple Normal Model

Suppose that each  $\theta_i \sim N(\mu, \tau^2)$ , that is, suppose each  $\theta_i$  is normally distributed with mean  $\mu$  and variance  $\tau^2$ . Suppose further that the  $\theta_i$  are independent. Let  $X_i$  given  $\theta_i$  be  $N(\theta_i, 1)$  and let the  $X_i$  be independent,  $i = 1, \dots, m$ . We shall use this simple model as a starting point.

For known  $\mu$  and  $\tau^2$ , the posterior distribution of  $\theta_i$  given  $\mathbf{X}$  is normal with mean  $E(\theta_i|\mathbf{X}) = \mu + \frac{\tau^2}{1+\tau^2}(X_i - \mu)$  and variance  $\text{var}(\theta_i|\mathbf{X}) = \frac{\tau^2}{1+\tau^2}$ . The  $\theta_i|\mathbf{X}$  are independent. Letting  $\gamma = \frac{\tau^2}{1+\tau^2}$ , we have

$$\begin{aligned} E\{G_m(t)|\mathbf{X}\} &= \frac{1}{m} \sum_{i=1}^m E\{I(\theta_i \leq t)|\mathbf{X}\} \\ &= \frac{1}{m} \sum_{i=1}^m \Pr(\theta_i \leq t|\mathbf{X}) \\ &= \frac{1}{m} \sum_{i=1}^m \Phi \left( \frac{t - \mu - \gamma(X_i - \mu)}{\sqrt{\gamma}} \right) \end{aligned} \quad (4.1)$$

where  $\Phi$  is the standard normal distribution function.

Let us consider the distance function

$$\|\hat{G}_m - G_m\|_{W,2} = \int_{-\infty}^{\infty} (\hat{G}_m(t) - G_m(t))^2 dW(t)$$

where  $W(t) \geq 0$  and  $\int_{-\infty}^{\infty} dW(t) < \infty$ . The conditional expected distance given  $\mathbf{X}$  is

$$\begin{aligned} E(\|\hat{G}_m - G_m\|_{W,2} | \mathbf{X}) &= E \left\{ \int_{-\infty}^{\infty} (\hat{G}_m(t) - G_m(t))^2 dW(t) \middle| \mathbf{X} \right\} \\ &= \int_{-\infty}^{\infty} E \left\{ (\hat{G}_m(t) - G_m(t))^2 \middle| \mathbf{X} \right\} dW(t). \end{aligned}$$

The last step is justified because the integrand is nonnegative and bounded. But the last integral can

be minimized by minimizing

$$E \left\{ (\hat{G}_m(t) - G_m(t))^2 \middle| \mathbf{X} \right\} \quad (4.2)$$

for each  $t$ . Note that the solution does not depend on  $W(t)$ . It is known from standard results in Bayes estimation that (4.2) is minimized by the choice  $\hat{G}_m(t) = E\{G_m(t)|\mathbf{X}\}$ . For the simple normal model, the latter quantity is given by (4.1).

*Note:* For  $W(t) \equiv t$ , Shen and Louis (1998) obtain

$$\hat{G}_m(t) = E\{G_m(t)|\mathbf{X}\} = \frac{1}{m} \sum_{i=1}^m \Pr(\theta_i \leq t|\mathbf{X})$$

for a two-stage hierarchical model.

It is of interest to compute the (conditional) expected loss because this provides a measure of the closeness of estimation, analogous to mean squared error. If  $\hat{G}_m(t) = E\{G_m(t)|\mathbf{X}\}$ , then

$$E \left\{ (\hat{G}_m(t) - G_m(t))^2 \middle| \mathbf{X} \right\} = \text{var}\{G_m(t)|\mathbf{X}\}, \text{ so}$$

$$E(\|\hat{G}_m - G_m\|_{W,2} | \mathbf{X}) = \int_{-\infty}^{\infty} \text{var}\{G_m(t)|\mathbf{X}\} dW(t). \quad (4.3)$$

But

$$\begin{aligned} \text{var}\{G_m(t)|\mathbf{X}\} &= \frac{1}{m^2} \sum_{i=1}^m \text{var}\{I(\theta_i \leq t)|\mathbf{X}\} \\ &= \frac{1}{m^2} \sum_{i=1}^m \Pr(\theta_i \leq t|\mathbf{X}) \{1 - \Pr(\theta_i \leq t|\mathbf{X})\} \\ &= \frac{1}{m^2} \sum_{i=1}^m \left[ \Phi \left( \frac{t - \mu - \gamma(X_i - \mu)}{\sqrt{\gamma}} \right) \right. \\ &\quad \left. \times \left\{ 1 - \Phi \left( \frac{t - \mu - \gamma(X_i - \mu)}{\sqrt{\gamma}} \right) \right\} \right]. \end{aligned} \quad (4.4)$$

From (4.3) and (4.4),  $E(\|\hat{G}_m - G_m\|_{W,2} | \mathbf{X})$  can be computed.

### 5. More General Situations

#### 5.1 Normal Model, Unequal but Known Variances

As before, let each  $\theta_i \sim N(\mu, \tau^2)$  and let the  $\theta_i$  be independent. Now let  $X_i$  given  $\theta_i$  be  $N(\theta_i, \sigma_i^2)$  and let the  $X_i$  be independent and  $\sigma_i^2 > 0$ ,  $i = 1, \dots, m$ .

Let

$$\gamma_i = \frac{\tau^2}{\sigma_i^2 + \tau^2}.$$

For known  $\mu$ ,  $\tau^2$ , and  $\sigma_i^2$ , the posterior distribution of  $\theta_i$  given  $\mathbf{X}$  is normal with mean

$$E(\theta_i|\mathbf{X}) = \mu + \gamma_i(X_i - \mu)$$

and variance

$$\text{var}(\theta_i|\mathbf{X}) = \gamma_i\sigma_i^2.$$

All the results of Section 4 continue to hold for this more general model, with

$$\frac{t - \mu - \gamma_i(X_i - \mu)}{\sigma_i\sqrt{\gamma_i}}$$

replacing

$$\frac{t - \mu - \gamma(X_i - \mu)}{\sqrt{\gamma}}$$

in (4.1) and (4.4).

## 5.2 Empirical and Hierarchical Bayes Techniques

Most frequently,  $\mu$  and  $\tau^2$  will be unknown and require estimation. There are standard empirical and hierarchical Bayes methods for doing this. See, for example, Ghosh and Rao (1994). The (conditional) expected loss can be estimated by means of Markov chain Monte Carlo methods.

## 6. Concluding Remarks

This paper has built upon the work of Louis (1984), Ghosh (1992), and others that study ways of estimating the distribution of small area parameters. Our focus has been on using loss functions that measure the distance between the distribution of the estimates of the parameters and the distribution of the parameters themselves. There are many aspects of this problem that have yet to be explored.

**Acknowledgment:** The author thanks Professor Malay Ghosh for acquainting him with this area of research in his 1994 Conference Board for the Mathematical Sciences lectures at the University of Connecticut, Storrs.

## REFERENCES

- Ghosh, M. (1992). Constrained Bayes estimation with applications, *Journal of the American Statistical Association* **87** 533–540.
- Ghosh, M. (1994). Bayesian methods in survey sampling, Conference Board for the Mathematical Sciences Seminar, University of Connecticut, unpublished notes.
- Ghosh, M., and Rao, J.N.K. (1994). Small area estimation: an appraisal (with discussion), *Statistical Science* **9** 55–93.
- Louis, T.A. (1984). Estimating a population of parameter values using Bayes and empirical Bayes methods, *Journal of the American Statistical Association* **79** 393–398.
- Rubin, D.B. (1981). Estimation in parallel randomized experiments, *Journal of Educational Statistics* **6** 377–401.
- Shen, W., and Louis, T.A. (1998). Triple-goal estimates in two-stage hierarchical models, *Journal of the Royal Statistical Society, Series B* **60** 455–471.

---

\*This paper is intended to promote the exchange of ideas among researchers and policy makers. The views are those of the author, and no official support by the U.S. Department of Education is intended or should be inferred.