DISSERTATION


PENALIZED ESTIMATION FOR SAMPLE SURVEYS IN THE PRESENCE OF

AUXILIARY VARIABLES


Submitted by

Mark J. Delorey

Department of Statistics


In partial fulfillment of the requirements

for the Degree of Doctor of Philosophy

Colorado State University

Fort Collins, Colorado

Summer 2008

UMI Number: 3332772

INFORMATION TO USERS

The quality of this reproduction is dependent upon the quality of the copy submitted. Broken or indistinct print, colored or poor quality illustrations and photographs, print bleed-through, substandard margins, and improper alignment can adversely affect reproduction.

In the unlikely event that the author did not send a complete manuscript and there are missing pages, these will be noted. Also, if unauthorized copyright material had to be removed, a note will indicate the deletion.

COLORADO STATE UNIVERSITY

May 30, 2008

WE HEREBY RECOMMEND THAT THE DISSERTATION PREPARED UNDER
OUR SUPERVISION BY MARK J. DELOREY ENTITLED PENALIZED ESTI-
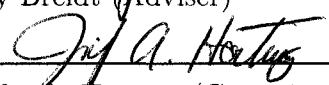MATION FOR SAMPLE SURVEYS IN THE PRESENCE OF AUXILIARY VARI-
ABLES BE ACCEPTED AS FULFILLING IN PART REQUIREMENTS FOR THE
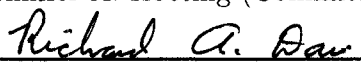DEGREE OF DOCTOR OF PHILOSOPHY.

Committee on Graduate Work

_____
F. Jay Breidt (Adviser)

_____
Jennifer A. Hoeting (Committee Member)

_____
Richard A. Davis (Committee Member)

_____
David M. Theobald (Committee Member)

_____
F. Jay Breidt (Department Head)

ABSTRACT OF DISSERTATION

PENALIZED ESTIMATION FOR SAMPLE SURVEYS IN THE PRESENCE OF
AUXILIARY VARIABLES

In conducting sample surveys, time and financial resources can be limited but
research questions are wide and varied. Thus, methods for analysis must make the
best use of whatever data are available and produce results that address a variety
of needs. Motivation for this research comes from surveys of aquatic resources, in
which sample sizes are small to moderate, but auxiliary information is available to
supplement measured survey responses. The problems of survey estimation are con-
sidered, tied together in their use of constrained/penalized estimation techniques for
combining information from the auxiliary information and the responses of interest.

We study a small area problem with the goal of obtaining a good ensemble es-
timate, that is, a collection of estimates for individual small areas that collectively
give a good estimate of the overall distribution function across small areas. Often,
estimators that are good for one purpose may not be good for others. For example,
estimation of the distribution function itself (as in Cordy and Thomas, 1997) can
address questions of variability and extremes but does not provide individual estima-
tors of the small areas, nor is it appropriate when auxiliary information can be made
of use. Bayes estimators are good individual estimators in terms of mean squared
error but are not variable enough to represent ensemble traits (Ghosh, 1992).

An algorithm that extends the constrained Bayes (CB) methods of Louis (1984)
and Ghosh (1992) for use in a model with a general covariance matrix is presented.
This algorithm produces estimators with similar properties as (CB), and we refer to

this method as general constrained Bayes (GCB). The ensemble GCB estimator is asymptotically unbiased for the posterior mean of the empirical distribution function (edf). The ensemble properties of transformed GCB estimates are investigated to determine if the desirable ensemble characteristics displayed by the GCB estimator are preserved under such transformations. The GCB algorithm is then applied to complex models such as conditional autoregressive spatial models and to penalized spline models. Illustrative examples include the estimation of lip cancer risk, mean water acidity, and rates of change in water acidity.

We also study a moderate area problem in which the goal is to derive a set of survey weights that can be applied to each study variable with reasonable predictive results. Zheng and Little (2003) use penalized spline regression in a model-based approach for finite population estimation in a two-stage sample when predictor variables are available. Breidt et al. (2005) propose a class of model-assisted estimators based on penalized spline regression in single stage sampling. Because unbiasedness of the model-based estimator requires that the model be correctly specified, we look at extending model-assisted estimation to the two-stage case. By calibrating the degrees of freedom of the smooth to the most important study variables, a set of weights can be obtained that produce design consistent estimators for all study variables. The model-assisted estimator is compared to other estimators in a simulation study. Results from the simulation study show that the model-assisted estimator is comparable to other estimators when the model is correctly specified and generally superior when the model is incorrectly specified.

Mark J. Delorey
Department of Statistics
Colorado State University
Fort Collins, Colorado 80523
Summer 2008

# ACKNOWLEDGEMENTS

# DEDICATION

This work is dedicated to Andrea and Jack.

CONTENTS

LIST OF FIGURES

LIST OF TABLES

# Chapter 1

# INTRODUCTION

A sample survey is generally undertaken to provide information about a population regarding particular research questions of interest. Often times, in addition to obtaining information about the whole population, the sample survey can be designed to provide information about subpopulations or *domains*, which may be defined by geographic areas or socio-demographic groups. Domains may be regarded as large or small, according to Rao (2003, Chapter 1); the definition depends upon whether or not the domain-specific sample is large enough to provide *direct estimates* of adequate precision. A *direct estimate* is based solely on the domain-specific sample data. If the domain-specific sample is large enough, the domain is regarded as "large"; otherwise, it is regarded as "small". Although it may use auxiliary information or be motivated by a model, a direct estimator is typically *design based* in the sense that inferences depend upon the probability distribution induced by the sampling design with the population values held fixed. In addition to the small and large domains, we introduce the notion of a *moderate* domain. For the purposes of this paper, we will distinguish a moderate domain as one in which direct estimates may be of adequate precision; however, the presence of auxiliary information suggests the estimate can be enhanced by using a model. This dissertation consists of two broad themes linked in their use of auxiliary data for estimation in sample surveys. The first theme involves the problem of small area estimation and the use of constrained Bayes methods in this context. Because our motivating examples are environmental, we are interested in extensions of these methods to spatial covariance models and to

semiparametric spatial regression models. Small area estimation is reviewed in Section 1.1. The second theme involves the problem of moderate area estimation using model-assisted survey regression estimates in the context of two-stages of unequal probability sampling. It is tied to the first research theme through consideration of penalized splines but is otherwise largely independent. Moderate area estimation is reviewed in Section 1.5.

## 1.1 Small Area Estimation

Because we will apply the small area estimation problem to natural resources data, including aquatic data, our domains of interest are geographic regions. If we were only interested in regional inferences, a probability sample over the entire region may be large enough to make model-free inferences using direct estimates. However, if the goal is to estimate responses for small geographic areas within this region, samples are not sufficiently dense in these small domains to make such estimates with reasonable precision. In fact, some of the small areas may contain no observations whatsoever. In this case, an *indirect* or *model-dependent* estimator is more appropriate. An indirect estimator will make use of responses in neighboring or related areas, effectively increasing the sample size in each of the small areas. The responses in neighboring areas are incorporated through a model which specifies the relationship between auxiliary data and the response as well as the relationship among responses in neighboring regions. Rao (2003, Chapter 5) classifies small area models into two types: *aggregate* or *area* level models that relate small area direct estimators to area-specific covariates, and *unit* level models that relate the individual units in the population to unit-specific covariates. We will focus on the area-level model.

For notation, let $\theta_h$ represent the characteristic of interest for small area $h$, where $h = 1, \ldots, m$. For example, in Section 3.3, $\theta_h$ represents acid neutralizing

capacity (a measure of water's acidity) and $h$ indexes the hydrological region as designated by the United States Geological Survey (Seaber et al., 1987). Also, let $\boldsymbol{x}_h = (x_{h1}, x_{h2}, \ldots, x_{hp})^T$ be a vector of known covariates associated with small area $h$, and let $\boldsymbol{\gamma} = (\gamma_1, \gamma_2, \ldots, \gamma_p)^T$ be a coefficient vector for the covariates. Then the basic area-level model, generally referred to as the *Fay-Herriot model* after Fay and Herriot (1979), is

$$\theta_h = \boldsymbol{x}_h^T \boldsymbol{\gamma} + \omega_h + \epsilon_h, \quad h = 1, \ldots, m. \tag{1.1}$$

Here, the $\omega_h$'s are area-specific random effects that are assumed to be independent and identically distributed (iid) with $E[\omega_h] = 0$ and $Var[\omega_h] = \sigma_\omega^2$, and the $\epsilon_h$'s are independent sampling errors with $E[\epsilon_h] = 0$ and $Var[\epsilon_h] = \psi_h^2$.

### 1.1.1 Small Area Estimation with Constrained Bayes

Inferential goals vary from study to study. In some cases, we are only interested in the individual area-specific characteristic, $\theta_h$. However, the variation in the characteristics among the small areas is also of interest in this situation. We therefore would like to find an ensemble estimate of the spatially-indexed true values $\{\theta_h\}_{h=1}^m$ that efficiently estimates the area-specific characteristics *and* whose distribution estimates the true empirical distribution function (edf)

$$F_m(z) = \frac{1}{m} \sum_{h=1}^m I\{\theta_h \le z\},$$

where $I\{A\} = 1$ if $A$ is true and $0$ otherwise. Following Rao (2003, p. 180), we extend the model specification to describe the uncertainty in estimating the spatially indexed $\{\theta_h\}_{h=1}^m$

$$y_h = \theta_h + \epsilon_h, \quad \{\epsilon_h\} \ NID\left(0, \ \sigma_h^2\right)$$

$$\theta_h = \boldsymbol{x}_h^T \boldsymbol{\gamma} + \omega_h, \quad \{\omega_h\} \ NID\left(0, \ \tau_\omega^2\right), \tag{1.2}$$

where *NID* denotes normally and independently distributed. We complete the probabilistic specification of all uncertainty in (1.1) by specifying a joint prior distribution

on $\sigma_h^2$, $\gamma$ and $\tau_\omega^2$ as

$$f\left(\sigma^2,\ \gamma,\ \tau_\omega^2\right) = g\left(\sigma^2\right) h\left(\gamma\right) k\left(\tau_\omega^2\right).$$

With this model, the posterior means

$$\theta_h^B = E\left(\theta_h \mid \boldsymbol{Y}\right) = E\left(\eta_h y_h + (1 - \eta_h)\, \mathbf{x}_h^T \boldsymbol{\gamma} \mid \boldsymbol{Y}\right), \qquad (1.3)$$

with

$$\eta_h = \frac{\tau_\omega^2}{\sigma_h^2 + \tau_\omega^2},$$

give the best mean squared error (MSE) *individual* estimates in the sense that

$$\sum_{h=1}^{m} E\left(\theta_h^B - \theta_h\right)^2 = \min_t \sum_{h=1}^{m} E\left(t_h - \theta_h\right)^2.$$

For fixed $\sigma_h^2$ and $\tau_\omega^2$, (1.3) can be written as

$$\eta_h y_h + (1 - \eta_h)\, x_h^T E\left(\gamma \mid \boldsymbol{Y}\right), \qquad (1.4)$$

in which it is easier to see that the $E\left(\theta_h \mid \boldsymbol{Y}\right)$ are weighted averages of the data (the $y_h$'s) and the mean predicted by the ($x_h^T \gamma$), given the data. If the variability of the prior distribution $\tau_\omega^2$ is large compared to the observation error $\sigma_h^2$, then $\eta_h$ will be large and the data is weighted heavier than the prior. The reverse is true when $\tau_\omega^2$ is small compared to $\sigma_h^2$.

However, the Bayes estimates in (1.3) are "over-shrunk"; there is too little variability among the $\left\{\theta_h^B\right\}$ to give a good representation of the edf of the $\theta_h$. Specifically, Louis (1984) and Ghosh (1992) show that

$$\sum_{h=1}^{m} \left(\theta_h^B - \bar{\theta}^B\right)^2 < E\left[\sum_{h=1}^{m} \left(\theta_h - \bar{\theta}\right)^2 \mid \boldsymbol{Y}\right],$$

where $\bar{\theta} = \frac{1}{m}\sum_{h=1}^{m}\theta_h$ and $\bar{\theta}^B = \frac{1}{m}\sum_{h=1}^{m}\theta_h^B$; that is, the empirical variability of the posterior means is strictly less than the posterior variability of the true values. Thus, the posterior means obtained from a Bayesian analysis are not good for *both*

individual and ensemble estimates. It is already the case that the sample mean of the Bayes estimates matches the posterior mean of $\bar{\theta}$, which is a desirable characteristic. Because the empirical variability of the Bayes estimates does not match the posterior variability of $\{\theta_h\}$, it is of interest for ensemble estimation to reduce the "shrinkage" so that the sample variance of the estimates matches the posterior variance of the true values (see Louis, 1984; Ghosh, 1992).

Following Ghosh (1992), first compute the scalars

$$H_1\left(\boldsymbol{Y}\right) = tr\left\{Var\left(\boldsymbol{\theta} - \bar{\theta}\boldsymbol{1} \mid \boldsymbol{Y}\right)\right\},$$

where $tr$ denotes the trace of a matrix, and

$$H_2\left(\boldsymbol{Y}\right) = \sum_{h=1}^{m}\left(\theta_h^B - \bar{\theta}^B\right)^2.$$

The *constrained Bayes* (CB) estimate of $\theta_h$ is then

$$\theta_h^{CB} = a\theta_h^B + \left(1 - a\right)\bar{\theta}^B \tag{1.5}$$

where

$$a = \left(1 + \frac{H_1\left(\boldsymbol{Y}\right)}{H_2\left(\boldsymbol{Y}\right)}\right)^{1/2} > 1.$$

Figure 1.1 shows alternative estimates of the edf of the $\{\theta_h\}$. Note the amount of shrinkage of the Bayes estimator. The posterior mean of $F_\theta$,

$$F_\theta^B\left(z\right) = \frac{1}{m}\sum_{h=1}^{m}E\left[I\left\{\theta_h \leq z\right\} \mid \boldsymbol{Y}\right],$$

is shown by the dashed curve and is the best estimate of the true edf in terms of MSE.

The posterior mean of $F_\theta$ does not, however, give an ensemble estimate in the sense of individual estimates of the $\{\theta_h\}$. The edf of the Bayes estimates and the edf of the CB estimates are both estimates of the true edf, and ensemble estimates. Figure 1.1 shows that the edf of the CB estimates is closer to the posterior mean of $F_\theta$

Figure 1.1: Shrinkage comparisons for ensemble estimates of $\{\theta_h\}$. The dashed line is the posterior mean of $F_\theta$, the gray line is $F_{CB}$, and the solid black line is $F_B$. The data are simulations from $y_h|\boldsymbol{\theta} \sim N(\theta_h, 7000)$, $\theta_h|\boldsymbol{\mu} \sim N(\mu_h, 4000)$, and $\mu_h \sim \text{uniform}(-100, 300)$ for $h = 1, \ldots, 200$.

than is the edf of the Bayes estimates. The edf of the Bayes estimates has too little variability (is "overshrunk") as an estimator of $F_\theta$, but the edf of the CB estimates is a better estimator of $F_\theta$ due to shrinkage reduction. This result is also shown numerically through a comparison of quantiles of the edf's to those of $F_\theta^B$. Table 1.1 shows that 0 is the $25^{th}$ percentile for the ensemble of Bayes estimators while it is at about the $31^{st}$ percentile for $F_\theta^B$ and for the constrained Bayes edf. These, combined with the estimates at $F_\theta(100)$, suggest that the edf of $\{\theta_h^{CB}\}$ is similar to $F_\theta^B$ and the edf of $\{\theta_h^B\}$ contains more mass at the center of the distribution.

As noted earlier, the small areas considered here are geographical regions. It might be expected that there exists some spatial correlation among the different regions. However, the form of CB estimates given in (1.5) does not take this into

Table 1.1: Comparison of ensemble estimates at selected quantiles. These numeric results are from the simulation presented graphically in Figure 1.1. The edf of $\{\theta_h^B\}$ contains almost 20% of its mass above 200 while the posterior mean of $F_\theta$ and the edf of $\{\theta_h^{CB}\}$ contain nearly 25% of their masses above 200. Likewise, the edf of $\{\theta_h^B\}$ contains 25% of its mass below 0 while the posterior mean of $F_\theta$ and the edf of $\{\theta_h^{CB}\}$ contain over 30% of their masses below 0.

| Estimate | $F_\theta(0)$ | $F_\theta(200)$ |
|---|---|---|
| edf of $\{\theta_h^B\}$ | 0.250 | 0.805 |
| posterior mean of $F_\theta$ | 0.307 | 0.766 |
| edf of $\{\theta_h^{CB}\}$ | 0.315 | 0.755 |

account. In Chapter 2, we will illustrate Stern and Cressie's (1999) extension of constrained Bayes to spatial data. We extend this further to take full account of parameter uncertainty, developing a methodology that applies to the broader setting of the general linear model with covariance matrix, known up to the values of a small number of parameters. This covers the case of linear mixed models, where the parameters are variance components, and so also applies to penalized splines. Additionally, it is often the case that when measuring geographical information, inferences are desired on a scale different from that upon which measurements were taken. So, we also will examine the problem of constrained Bayes estimates on different scales.

## 1.2 Conditional auto-regressive models

In this section we review certain spatial models that are relevant to the problem of geographic small area estimation. Suppose $\{\theta_h\}_{h=1}^m$ follows a joint Gaussian distribution with means $\{\mu_h\}$. Let $\boldsymbol{\theta}_{-h} = (\theta_1, \ldots, \theta_{h-1}, \theta_{h+1}, \ldots, \theta_m)^T$. Then the

conditional density of $\theta_h$ given $\boldsymbol{\theta}_{-h}$ can be written as

$$p\left(\theta_h|\boldsymbol{\theta}_{-h}\right) = \frac{1}{\sqrt{2\pi\left(\sigma v_h\right)^2}}\exp\left\{-\frac{1}{2\left(\sigma v_h\right)^2}\left[\theta_h - \eta_h\left(\boldsymbol{\theta}_{-h}\right)\right]^2\right\}, \qquad (1.6)$$

where $\eta_h\left(\boldsymbol{\theta}_{-h}\right)$ is the conditional mean and $\left(\sigma v_h\right)^2$ is the conditional variance. Besag (1974) and Cressie (1993) show that, under a regularity condition of pair-wise only dependence between lattice points,

$$\eta_h\left(\boldsymbol{\theta}_{-h}\right) = \mu_h + \sum_{\substack{i=1 \\ i\neq h}}^{m}c_{hi}\left(\theta_h - \mu_i\right), \quad h = 1,\ldots,m, \qquad (1.7)$$

where $\mu_h$ is the unconditional mean of $\theta_h$, $c_{hi}v_i^2 = c_{ih}v_h^2$, $c_{hh} = 0$, and $c_{hk} = 0$ unless there is pairwise dependence between area $h$ and area $k$. Additionally, if $\left(\boldsymbol{I} - \boldsymbol{C}\right)^{-1}$ exists and $\left(\boldsymbol{I} - \boldsymbol{C}\right)^{-1}\boldsymbol{V}$ is symmetric, where $\boldsymbol{I}$ is the identity matrix, $\boldsymbol{C} = \left(c_{hk}\right)$, and $\boldsymbol{V} = \text{diag}\left(v_1,\ldots,v_m\right)$, then

$$\boldsymbol{\theta} \sim N\left(\boldsymbol{\mu}, \sigma^2\left(\boldsymbol{I} - \boldsymbol{C}\right)^{-1}\boldsymbol{V}\right). \qquad (1.8)$$

Cressie (1993) notes that model (1.7) qualifies as a spatial analogue of an autoregressive time series since, if we define $\boldsymbol{\omega} \equiv \left(\boldsymbol{I} - \boldsymbol{C}\right)\left(\boldsymbol{\theta} - \boldsymbol{\mu}\right)$, then

$$\theta_h - \mu_h \equiv \sum_{\substack{i=1 \\ i\neq h}}^{m}c_{hi}\left(\theta_h - \mu_i\right) + \omega_h, \quad h = 1,\ldots,m, \qquad (1.9)$$

which has a form similar to that of the AR($p$) time series

$$X_t = \sum_{i=1}^{p}\phi_i X_{t-i} + Z_t,$$

where $\{Z_t\}$ is a zero-mean process with $Z_s$ independent of $X_t$ for all $s < t$. Likewise, under the definition in (1.9), $E\left[\boldsymbol{\omega}\right] = 0$ and

$$
\begin{aligned}
Cov\left[\boldsymbol{\omega}, \boldsymbol{\theta}\right] &= E\left[\boldsymbol{\omega}\boldsymbol{\theta}^T\right] \\
&= E\left[\boldsymbol{\theta}\boldsymbol{\theta}^T - \boldsymbol{\mu}\boldsymbol{\theta}^T - \boldsymbol{C}\boldsymbol{\theta}\boldsymbol{\theta}^T + \boldsymbol{C}\boldsymbol{\mu}\boldsymbol{\theta}^T\right] \\
&= \sigma^2\left(\boldsymbol{I} - \boldsymbol{C}\right)^{-1}\boldsymbol{V} + \boldsymbol{\mu}\boldsymbol{\mu}^T - \boldsymbol{\mu}\boldsymbol{\mu}^T - \sigma^2\boldsymbol{C}\left(\boldsymbol{I} - \boldsymbol{C}\right)^{-1}\boldsymbol{V} - \boldsymbol{C}\boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{C}\boldsymbol{\mu}\boldsymbol{\mu}^T \\
&= \sigma^2\left(\boldsymbol{I} - \boldsymbol{C}\right)\left(\boldsymbol{I} - \boldsymbol{C}\right)^{-1}\boldsymbol{V} \\
&= \sigma^2\boldsymbol{V},
\end{aligned}
$$

which is a diagonal matrix. This implies, due to normality, that $\omega_h$ is independent of $\theta_{-h}$ for all $h$. Because of the properties outlined here, the model (1.8) is called a *conditional autoregressive* or CAR model.

Banerjee et al. (2004) note that since the conditional densities given in (1.7) are compatible, an analytical expression for the joint distribution of $\theta$ given by (1.8) is ensured by Brook's lemma. The conditions $c_{hi}v_i^2 = c_{ih}v_h^2$ are still needed to guarantee that $(I - C)^{-1}V$ is symmetric. However, they give an example of a case in which the conditionals (1.7) are defined, yet the covariance matrix $(I - C)^{-1}V$ in (1.8) is singular, and thus, the distribution of $\theta$ is improper. They caution, then, that care must be taken in defining the matrices $C$ and $V$. Both Banerjee et al. (2004) and Cressie (1993) show that the impropriety problem can be taken care of by redefining (1.8) as

$$\theta \sim N\left(\mu, \sigma^2 \left(I - \phi C\right)^{-1} V\right) \tag{1.10}$$

and choosing $\phi$ so that $(I - \phi C)^{-1}V$ is positive-definite. Positive-definiteness is achieved if and only if

$$\phi \in (\phi_{\min}, \phi_{\max}), \tag{1.11}$$

where $\phi_{\min} = \eta_1^{-1}$, $\phi_{\max} = \eta_n^{-1}$, and $\eta_1 < 0 < \eta_n$ are the smallest and largest eigenvalues, respectively, of $V^{1/2}CV^{1/2}$.

The choices made for $\Phi = (I - \phi C)^{-1}V$ determine the spatial relationships. The matrix $C$ is called the *adjacency matrix* since its elements $c_{ij}$ can be interpreted as a measure of spatial association between area $i$ and area $j$. If $c_{ij} = 0$ then area $i$ and area $j$ are statistically independent. The larger the value of $c_{ij}$ the stronger the dependence between areas $i$ and $j$. This property can be used to reflect, for example, the hierarchical nesting of areas and subareas. The diagonal elements $v_{hh}$ of the matrix $V$ represent variance scaling factors for each $\theta_h$. The parameter $\phi$ is called the *spatial dependence parameter* and is a measure of the overall strength of spatial correlation. Clearly if $\phi = 0$, then $\Phi = \sigma^2 I$ and $\theta_h$ and $\theta_k$ are independent for all $h, k$.

## 1.3   Penalized splines

The general problem reviewed in this section involves the common action of fitting a model to paired data $(x_i, y_i)$, where $y_1, \ldots, y_n$ are responses that have been observed at fixed, non-stochastic design points $x_1 < \cdots < x_n$ and generated from the model

$$y_i = f(x_i) + \epsilon_i, \qquad i = 1, \ldots, n. \tag{1.12}$$

Here, $f$ is an unknown regression function and $\epsilon_1, \ldots, \epsilon_n$ are zero-mean, uncorrelated errors. The problem then is to estimate $f$.

A traditional method for the estimation of $f$ is to assume some parametric form for $f$ and minimize the residual sum of squares criterion

$$RSS\left(\tilde{f}\right) = \sum_{i=1}^{n} \left(y_i - \tilde{f}(x_i)\right)^2 \tag{1.13}$$

over all functions $\tilde{f}$ with the assumed parametric form. A problem with this approach arises when the function $f$ is not from the parametric family that is assumed. Ruppert et al. (2003) use, as an example, LIDAR (light detection and ranging) data shown in Figure 1.2 along with three polynomial fits. The explanatory variable "range" is the distance traveled before the light is reflected back to its source, and the response is the log of the ratio of the received light from two laser sources. Interesting features of these data include a distinct non-linear trend and heteroscedasticity (the response is much more variable for higher ranges than for lower ranges). As can be seen in Figure 1.2, neither the third nor fourth degree polynomial is flexible enough to follow the trend; the third degree polynomial does not capture the relatively flat region of the data over the lower ranges, while the fourth degree polynomial curves up unnecessarily at the end. The tenth degree polynomial does a better job, but has too many slope changes in the flat region of the lower ranges. Thus, even higher degree polynomials will not always provide an adequate fit to the data nor are they

are called *knots*, referring to the fact that the piecewise function is "tied together" at these points. The part of (1.14) involving the terms $\beta_0, \ldots, \beta_p$ is just a regular polynomial function. However, by including the knot terms, we are allowing the coefficients of $x^1, \ldots, x^p$ to change between each pair of knots. This can be seen by rewriting (1.14) as follows:

$$
\begin{aligned}
\tilde{f}(x_i) &= \beta_0 + \beta_1 x_i + \cdots + \beta_p x_i^p + \sum_{k=1}^{K} \beta_{p+k} \left[ (x_i - \kappa_k)_+ \right]^p \\
&= \beta_0 + \beta_1 x_i + \cdots + \beta_p x_i^p + \sum_{k=1}^{K} \beta_{p+k} \sum_{j=0}^{p} \binom{p}{j} x_i^j (-\kappa_k)^{p-j} I_{\{x_i - \kappa_k > 0\}} \\
&= \beta_0 + \beta_1 x_i + \cdots + \beta_p x_i^p \\
&\quad + \sum_{k=1}^{K} \beta_{p+k} I_{\{x_i - \kappa_k > 0\}} \left[ (-\kappa_k)^p + \binom{p}{1} x_i (-\kappa_k)^{p-1} + \cdots + x_i^p \right] \\
&= \left( \beta_0 + \sum_{k=1}^{K} \beta_{p+k} I_{\{x_i - \kappa_k > 0\}} (-\kappa_k)^p \right) \\
&\quad + \left( \beta_1 + \sum_{k=1}^{K} \beta_{p+k} I_{\{x_i - \kappa_k > 0\}} \binom{p}{1} (-\kappa_k)^{p-1} \right) x_i \\
&\quad + \cdots + \left( \beta_0 + \sum_{k=1}^{K} \beta_{p+k} I_{\{x_i - \kappa_k > 0\}} \right) x_i^p.
\end{aligned}
\tag{1.15}
$$

If we say $f_k(x) = \tilde{f}(x)$ for $\kappa_k < x \le \kappa_{k+1}$, equation (1.15) shows that, for each $k$, $f_k$ is a (potentially) different polynomial of degree $p$.

The spline in (1.14) uses knot terms of the form $(x - \kappa_k)_+^p$. Note that any function that is piecewise $p$-th order polynomial, with knots at $\kappa_1, \ldots, \kappa_K$, can be obtained as a linear combination of the functions $\left\{ 1, x, \ldots, x^p, (x - \kappa_1)_+^p, \ldots, (x - \kappa_K)_+^p \right\}$. Thus, we call this set a *spline basis*, more specifically a *truncated power basis of degree p*. This is by no means the only spline basis that we could have chosen.

One of the disadvantages of the truncated power basis is that the functions are not orthogonal. Computationally, this can lead to numerical instability when there are a large number of knots or the smoothing parameter is close to zero. B-splines

(see Eilers and Marx, 1996) are the result of a transformation of the truncated power basis and lead to more stable estimation of the parameters. Other bases include the trigonometric basis, Demmler-Reinsch basis, and wavelet basis (see Nychka and Cummins, 1996; Ogden, 1996; Hardle et al., 1998).

Because of their interpretability, we will assume the use of the truncated power basis functions. In fitting a spline, we would like to be able to capture the trend in the scatter without picking up microvariation or overfitting to chance fluctuations. The greater the number of knots in the model, the more frequently we are allowing $\hat{f}$ to change its slope, and thus the closer we come to interpolating the data. For example, assuming we have only one response for each $x_i$, if we let the knots occur at the design points $x_1, \ldots, x_n$, then $\tilde{f}$ will be a $p$-th order piecewise polynomial that interpolates the data, and hence will be severely overfitted. Two ways to reduce this overfitting are either to reduce the number of knots or to constrain the influence of the knots themselves. Ruppert et al. (2003, p. 65) suggest that the constraint $\sum \beta_k^2 < C$ is one that can rectify overfitting and is easy to implement. If we let $\boldsymbol{\beta} = (\beta_0, \ldots, \beta_p, \beta_{p+1}, \ldots, \beta_K)^T$, where $\beta_0, \ldots, \beta_p$ are the coefficients of $x^0, \ldots, x^p$, respectively, and $\beta_{p+1}, \ldots, \beta_K$ are the coefficients of the knot terms, then $\sum_{k=p+1}^{K} \beta_k^2 < C$ can be written as $\boldsymbol{\beta}^T \mathbf{D} \boldsymbol{\beta} < C$, where

$$D = \begin{bmatrix} \mathbf{0}_{(p+1)\times(p+1)} & \mathbf{0}_{(p+1)\times K} \\ \mathbf{0}_{K\times(p+1)} & \boldsymbol{I}_{K\times K} \end{bmatrix}.$$

The problem then becomes one in which we want to minimize (1.13) subject to $\boldsymbol{\beta}^T \mathbf{D} \boldsymbol{\beta} < C$. By using a Lagrange multiplier argument, this means choosing $\boldsymbol{\beta}$ to minimize

$$\|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|^2 + \lambda \boldsymbol{\beta}^T \mathbf{D} \boldsymbol{\beta}. \tag{1.16}$$

This minimization problem has the solution

$$\hat{\boldsymbol{\beta}} = \left(\boldsymbol{X}^T \boldsymbol{X} + \lambda \boldsymbol{D}\right)^{-1} \boldsymbol{X}^T \boldsymbol{y}.$$

The larger the sum $\sum_{k=p+1}^{K} \beta_k^2$, the more significant is the change of slope at each of the knot points, and thus, the "rougher" the estimated curve will be. Hence, $\lambda \beta^T \mathbf{D} \beta$ is called the *roughness penalty* since, for fits that are rough, this term will be large, and thus, these fits are unlikely to be selected. The parameter $\lambda$ is called the *smoothness parameter* because it determines the weight or importance that is placed on the smoothness of the fit. If $\lambda = 0$, then $\hat{f}$, in the extreme case where there are knots at each design point, is just the interpolant. As $\lambda \rightarrow \infty$, a change of slope at the knots is not allowed, and $\hat{f}$ approaches a polynomial of degree $p$.

### 1.3.1    Selecting knot locations and the value of the smoothing parameter

It was noted earlier that the number and location of the knots, together with the value of the smoothing parameter, $\lambda$, affect how well the fitted function can respond to changes in variability of the data as well as how smooth the fit will be. It is desirable to place more knots where the response is more highly variable; doing so will allow the rate at which $\hat{f}$ changes slope to be higher where needed. Conversely, fewer knots are needed where the response is less variable. When the knots are at the design points $x_1, \ldots, x_n$, the result is a *smoothing spline*. Green and Silverman (1994) note that, though the use of the design points as knots leads to an exact fit, it may be desirable to approximate this fit with a smaller number of knots. Reasons for doing so include large sample sizes that make the solution computationally intense and when a smaller number of parameters are desired for ease of interpretation. Ruppert et al. (2003) give details on how so called *low rank smoothers*, smoothers that use significantly fewer knots than data points, extract only the essential information from the data. There are several data driven methods for estimating the smoothing parameter: likelihood approach, cross-validation (CV) and generalized cross-validation (GVC), Mallows $C_p$, and AICC, which are all described by Ruppert et al. (2003).

## 1.3.2 Degrees of freedom of the smooth

Ruppert et al. (2003) describe how it is possible to measure the degrees of freedom of a smooth. This will assist us later on in interpreting the smooth since a penalized spline with $\nu$ degrees of freedom smooths the trend about the same amount, in some sense, as a $\nu$th degree polynomial. In parametric regression, the *hat matrix* $\boldsymbol{H}$ is so called because it converts the observed response $\boldsymbol{y}$ into the predicted response $\hat{\boldsymbol{y}}$, i.e.,

$$\hat{\boldsymbol{y}} = \boldsymbol{H}\boldsymbol{y}. \tag{1.17}$$

It is a well known result of simple linear modeling (see, for example, Hocking, 1996, p. 302) that

$$
\begin{aligned}
tr\left(\boldsymbol{H}\right) &= \text{ number of fitted parameters}\\
&= \text{ degrees of freedom.}
\end{aligned}
\tag{1.18}
$$

The model in (1.14) can be written as

$$\tilde{f}\left(x\right) = \boldsymbol{X}\boldsymbol{\beta}$$

and estimated by

$$
\begin{aligned}
\hat{f}\left(x\right) &= \boldsymbol{X}\hat{\boldsymbol{\beta}}\\
&= \boldsymbol{X}\left(\boldsymbol{X}^T\boldsymbol{X} + \lambda\boldsymbol{D}\right)^{-1}\boldsymbol{X}^T\boldsymbol{y}\\
&= \boldsymbol{S}_\lambda\boldsymbol{y}.
\end{aligned}
$$

Here, $\boldsymbol{S}_\lambda$, called the *smoother matrix*, plays a role similar to $\boldsymbol{H}$ in equation (1.17). Both are linear operators on $\boldsymbol{y}$ and convert an observed response vector to a vector of predicted responses. Extending the definition of *degrees of freedom* given in (1.18), we can say that

$$
\begin{aligned}
tr\left(\boldsymbol{S}_\lambda\right) &= \text{ equivalent number of parameters}\\
&= \text{ degrees of freedom of the smooth.}
\end{aligned}
$$

As shown by Ruppert et al. (2003), if a penalized spline has $\kappa$ knots and degree $p$, then

$$tr\left(\boldsymbol{S}_0\right) = p + 1 + \kappa$$

and

$$tr\left(\boldsymbol{S}_\lambda\right) = p + 1 \quad \text{as } \lambda \to \infty,$$

which corresponds to the piecewise polynomial and the global polynomial, respectively.

### 1.3.3 Extension to spatial regression

In the introduction to Section 1.3, we note that we use the truncated power basis functions because of their interpretability. This is true for the model-assisted estimation in Chapter 5 in which the data are one-dimensional. In Chapter 4 we apply a general constrained Bayes algorithm to spatial data that are modeled using a penalized spline. The model in (1.12) is extended to the two-dimensional model

$$y_i = f\left(s_i, t_i\right) + \epsilon_i, \qquad i = 1, \ldots, n, \tag{1.19}$$

where $f$ is a real-valued bivariate function. Once again, $\tilde{f}$ is an estimator of $f$, and we restrict $\tilde{f}$ to the class of penalized splines. The spline model in (1.14) can be naturally extended using a two-dimensional version of the truncated power basis. For example, with $p = 1$, we have (Ruppert et al., 2003, p. 240)

$$
\begin{aligned}
\tilde{f}\left(x_i\right) &= \beta_0 + \beta_s s_i + \beta_t t_i + \sum_{k=1}^{K_s} u_k^s \left(s_i - \kappa_k^s\right)_+ + \sum_{k=1}^{K_t} u_k^t \left(t_i - \kappa_k^t\right)_+ \\
&\quad + \gamma s_i t_i + \sum_{k=1}^{K_s} v_k^s s_i \left(t_i - \kappa_k^s\right)_+ + \sum_{k=1}^{K_t} v_k^t t_i \left(s_i - \kappa_k^t\right)_+ \\
&\quad + \sum_{k=1}^{K_s} \sum_{k=1}^{K_t} v_{kk'}^{st} \left(s_i - \kappa_k^s\right)_+ \left(t_i - \kappa_k^t\right)_+ .
\end{aligned}
\tag{1.20}
$$

Because the model in (1.20) is obtained by forming all pairwise products of the basis functions

$$1, s, \left(s - \kappa_1^s\right)_+, \ldots, \left(s - \kappa_{K_s}^s\right)_+$$
$$1, t, \left(t - \kappa_1^t\right)_+, \ldots, \left(t - \kappa_{K_t}^t\right)_+,$$

it is referred to as the *tensor product* basis. However, the tensor product basis is not rotational invariant (Ruppert et al., 2003). Therefore, the analysis of any data dependent upon the orientation of their coordinate system (such as geographical data) would yield different results if data are measured on a different set of axes. Ruppert et al. (2003) then note that rotational invariance can be achieved by using a *radial basis*. These are functions of the form

$$C\left(\left\| (s, t) - \left(\kappa^s, \kappa^t\right) \right\|\right), \tag{1.21}$$

where $C$ is a univariate function. Since $C(s, t)$ depends only on the distance, and not the direction, from the knot $(\kappa^s, \kappa^t)$, reorienting the coordinate system will not affect the value of $C$ at $(s, t)$.

## 1.4 Lattice models and penalized splines as general linear models

The simple linear model

$$\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \text{ where } \boldsymbol{\epsilon} \sim N\left(\boldsymbol{0}, \sigma^2 \boldsymbol{I}\right) \tag{1.22}$$

assumes heteroscedasticity and independence among the $y_i$'s, among other assumptions. When a more complex variance structure exists in a linear model, one approach is to generalize (1.22) to

$$\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \text{ where } \boldsymbol{\epsilon} \sim N\left[\boldsymbol{0}, \sigma^2 \boldsymbol{V}\left(\boldsymbol{\phi}\right)\right] \tag{1.23}$$

and where $\boldsymbol{V}\left(\boldsymbol{\phi}\right)$ is a positive definite matrix, possibly depending upon some vector of parameters $\boldsymbol{\phi}$. The linear mixed model

$$\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{Z}\boldsymbol{u} + \boldsymbol{\epsilon}, \tag{1.24}$$

where $u \sim N(0, G)$ and $\epsilon$ is as in (1.22) can be written in the form of (1.23) by defining $\epsilon* = Zu + \epsilon$ and calling $V = V(\phi) = ZGZ' + \sigma^2 I$. The *generalized least squares* or gls estimate of $\beta$ under model (1.23) is $\hat{\beta} = (X'V^{-1}X)^{-1} X'V^{-1}Y$. Furthermore, if $V(\phi)$ is known, this is the *best linear unbiased estimator* (BLUE) of $\beta$ McCulloch and Searle (see, for example 2001). There are many models that fit into this framework. For example, the first-order moving average process

$$y_t = Z_t + \alpha Z_{t-1}$$
$$\{Z_t\}_{t=1}^n \sim N(0, \sigma^2 I)$$

can be written as

$$y \sim N(0, \sigma^2 V),$$

where

$$V = \begin{bmatrix} \alpha^2+1 & \alpha & 0 & 0 & 0 & \cdots & 0 \\ \alpha & \alpha^2+1 & \alpha & 0 & 0 & \cdots & 0 \\ 0 & \alpha & \alpha^2+1 & \alpha & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \alpha & \alpha^2+1 & \alpha \\ 0 & 0 & 0 & 0 & 0 & \alpha & \alpha^2+1 \end{bmatrix}.$$

It is easy to see that the CAR spatial model (1.8) is of the form in (1.23). Stern and Cressie (1999) make use of this fact in deriving CB estimates for the CAR model, and we make use of it to extend their method to apply to more general problems. Model (1.12) with $f$ a penalized spline is also of the form in (1.23) which we demonstrate here by outlining Ruppert et al. (2003), Chapters 3 and 4. Recall the penalized spline fitting criterion $\|y - X\beta\|^2 + \lambda\beta^T D\beta = \|y - X\beta\|^2 + \lambda \|\beta_2\|^2$ in (1.16), where $\beta_2 = (\beta_{p+1}, \ldots, \beta_K)'$. Dividing this by $\sigma_\epsilon^2$ yields

$$\frac{1}{\sigma_\epsilon^2} \|y - X\beta\|^2 + \frac{\lambda}{\sigma_\epsilon^2} \|\beta_2\|^2. \tag{1.25}$$

The expression in (1.25) can be made equal to the BLUE criterion by considering $\beta_2$ a set of random coefficients with $Var(\beta_2) = \sigma_{\beta_2}^2 I$, where $\sigma_{\beta_2}^2 = \sigma_\epsilon^2/\lambda$. A standard

method for estimating the smoothing parameter $\lambda$ is to treat $\beta_{p+1}, \ldots, \beta_K$ in (1.14) as random effects in a linear mixed model specification. The result is that (1.12), with $f$ as in (1.14) can be written as

$$Y = X\beta_1 + Z\beta_2 + \epsilon, \tag{1.26}$$

where $\beta_1 = (\beta_0, \ldots, \beta_p)'$, $\beta_2 = (\beta_{p+1}, \ldots, \beta_K)'$, and $X$ and $Z$ are known design matrices for the fixed effects $\beta_1$ and random effects $\beta_2$, respectively. Additionally, we let $\beta_2 \sim (0, \eta^2 I)$ and $\epsilon \sim (0, \sigma^2 I)$. Thus, (1.26) can be written as (1.23) with $V = \eta^2 ZZ' + \sigma^2 I$. In this formulation, the GLS estimator for $\beta_1$ is

$$\hat{\beta}_1 = \left(X'V^{-1}X\right)^{-1} X'V^{-1}Y. \tag{1.27}$$

## 1.5 Moderate area estimation

A *moderate* domain, as defined in the introduction, is one in which direct estimates may be of adequate precision, but the presence of auxiliary variables suggests enhancing the estimate by using a model. A model can improve upon the estimates, in terms of their variability, in the large area estimation problem, too. However, whenever a model is introduced, there is the risk that it may be misspecified. A grossly misspecified model can result in estimates for which the bias is quite large. In a true large area estimation problem, we may not want to risk misspecifying the model. The direct estimates have sufficient precision, and the additional gain in precision that can come through a model is not worth the risk of the increased bias if the model is misspecified. In a moderate area estimation problem, we may be content with the precision of the direct estimates. However, the improvement in precision obtained from introducing a model is worth the risk of misspecification. To alleviate some of the risk of misspecifying the model, we will use penalized splines, a semi-parametric alternative to the specification of a parametric model.

## 1.6   Structure of paper

Chapter 2 begins with a discussion of the method given by Stern and Cressie (1999) that incorporates spatial relationships into the constrained Bayes approach through a CAR model. As noted above, a CAR fits into the general linear model framework, and the CB estimators make use of the projection matrix derived from a gls regression. In Section 2.3 we extend the method of Stern and Cressie (1999) to a general covariance matrix and introduce a completely hierarchical Bayesian model. An algorithm for finding CB estimates for a linear model with a general covariance matrix is presented. Our applications in later chapters are ecological, in which there is often a need for estimators/predictors for different scales of geography. So, in the remaining sections of Chapter 2, we will examine the properties of constrained Bayes estimates on scales different from that of the data. Primarily, this involves studying properties of linear combinations of the estimates.

In Chapter 3 we compare the analysis of the Scotland lip cancer data done by Stern and Cressie (1999) to an analysis using our extension of their method. Both analyses use CB estimates derived from a CAR model. In Stern and Cressie (1999), the spatial association parameter $\phi$ is fixed, and they select the estimate of $\phi$ by running the analysis several times using different values of $\phi$ and choosing the one that yields the lowest mean square prediction error (MSPE). Our methods assume $\phi$ is random, and its value is estimated from the data. Continuing with a CAR model, we analyze the acid neutralizing capacity (ANC) for watersheds in the Mid-Atlantic Highlands of the United States. The watersheds have a nested structure to them, with watersheds within a particular nesting level sharing similar characteristics in terms of geography. Therefore, we can try to identify correlation among different nesting levels. We analyze the data using first one, then two levels of nesting.

In Chapter 4, we apply the algorithm developed in Chapter 2 to a penalized spline model. As noted in Section 1.4, Ruppert et al. (2003) demonstrate that

penalized splines can fit into the linear mixed model framework. Thus, our algorithm can be used to find CB estimates of the fixed and random components of the penalized spline model. We use the model and data on Northeast lakes from Opsomer et al. (2008)

Chapter 5 turns to the second research theme of this dissertation, with an exploration of moderate area estimation with penalized splines. We construct a model-assisted survey regression estimator for two-stage sampling using penalized splines, extending the single-stage estimator in Breidt et al. (2005). It is shown that the estimator is asymptotically design-unbiased and design consistent under mild assumptions and that its variance can be consistently estimated. The two-stage model-assisted estimator is contrasted with the estimator in Zheng and Little (2003), which uses penalized spline regression in a model-based approach. A series of simulations demonstrate that the model-assisted estimators generally fare no worse than model-based estimators when the model is correctly specified and generally are superior to model-based estimators when the model is incorrectly specified.

# Chapter 2

# CONSTRAINED BAYES ESTIMATION

## 2.1   Introduction

We approach the problem of small area estimation for an ensemble of characteristics from a Bayes perspective. Keeping in mind our two-fold objective, the limitations of the typical posterior mean have already been examined: while the posterior means are good for *individual* estimates, constrained Bayes estimation is more appropriate for the *ensemble* of estimates since the set of constrained Bayes estimators $\{\theta_h^{CB}\}$ matches the posterior mean and posterior variance of the true values $\{\theta_h\}_{h=1}^m$ (Ghosh, 1992). Stern and Cressie (1999) illustrate how constrained Bayes can be employed in a CAR model. We extend this idea to a fully Bayesian context with a general covariance matrix.

The mapping of disease incidence and mortality rates has been used by epidemiologists and others to identify so called "hotspots", geographic regions with unusually high (or low) rates. Not only are these extremes properly regarded as important components to the summary of the data, but extremes are also of interest as they are often indicative of unusual or extraordinary conditions. In terms of public health, extremely high or low disease rates are usually cause for intervention or special study. Various methods have been used to estimate these rates. Clayton and Kaldor (1987), Cressie and Chan (1989), and Kleinschmidt et al. (2001) employ an empirical Bayes approach using a conditional autoregressive model (CAR) that accounts for spatial relationships among neighboring regions. Manton et al. (1989)

also use empirical Bayes methods while accounting for spatial dependence. While these methods give good estimates of individual rates, they may not be appropriate for broader inferential goals. It has been pointed out by several authors (see, for example Louis, 1984; Ghosh, 1992; Devine et al., 1994a; Devine and Louis, 1994; Devine et al., 1994b; Stern and Cressie, 1999) that the posterior means used in standard Bayesian analysis are "overshrunk" and thus do not give a good representation of the overall distribution of these rates. Constrained Bayes methods Ghosh (1992) described in Section 1.1.1 force the variability of the estimates to match the posterior variability of the true values, thus providing an ensemble of estimates that satisfy certain first and second-order moment conditions and, therefore, reflect properties of the true vector of parameters.

## 2.2 Background

We now extend the model in (1.10) and let

$$Y \mid \theta, \sigma \sim N\left(\theta, \sigma^2 D\right)$$
$$\theta \mid \gamma, \tau_\omega, \phi \sim N\left(X\gamma, \tau_\omega^2 \left(I - \phi C\right)^{-1} V\right) \tag{2.1}$$

where $\sigma$ is an unknown parameter, $D$ is a known matrix (not necessarily diagonal), $X$ is a design matrix of covariates associated with each lattice point, and $\gamma$ is an unknown coefficient vector.

With this model specification, Bayes estimates of $\theta$ can be obtained relatively easily either analytically or numerically. A benefit of using Bayes estimates is that the impact of sparsely-sampled, high-variance areas is lessened by borrowing strength across areas through the model (2.1). Thus, using Bayes estimates seems to make sense. However, for the two-fold inference goal described in Chapter 1, in which we are interested in good *individual* estimates and a good *ensemble* estimate, the same limitations of the Bayes estimates arise. The Bayes estimates are "overshrunk"

toward the overall mean, and Gelman and Price (1999) show that, in spatial mapping, sparse areas are much less likely to appear as extreme values than dense areas. Therefore, if we are interested in capturing the true variability of the response over the lattice or identifying high and low extremes, Bayes estimates may not be the best choice.

As was discussed in Chapter 1, constrained Bayes (CB) estimates address this issue by forcing the estimates to match the posterior expected values of the sample mean and sample variance of the parameters. Stern and Cressie (1999) provide an extension of CB estimates to the spatial case with a CAR model. The main result we will use is Theorem 1 from Stern and Cressie (1999, p. 77). We state it here without proof and refer the reader to the paper for details of the proof.

**Theorem 1** *Suppose that $Y|\theta \sim p(Y|\theta)$ and $\theta \sim N(X\gamma, \sigma^2\Phi)$, with $\Phi$ positive definite, and $\gamma$, $\sigma$, and $\Phi$ known. Let $\hat{\theta}_0 = \{Y : H_2(Y) > 0\}$ with $H_2(Y)$ defined below, and let $P = X\left(X^T\Phi^{-1}X\right)^{-1}X^T\Phi^{-1}$ denote the projection matrix that yields the predicted values for the generalized least squares regression on $X$ with error vector that has variance matrix $\Phi$. Then, for $Y \in \hat{\theta}_0$, the estimator $t(Y)$ that minimizes the posterior expected weighted squared error $E\left[(\theta - t(Y))^T\Phi^{-1}(\theta - t(Y))|Y\right]$ subject to*

$$PE[\theta|Y] = Pt(Y) \qquad (2.2)$$

*and*

$$E\left[\theta^T(I-P)^T\Phi^{-1}(I-P)\theta|Y\right] = t(Y)^T(I-P)^T\Phi^{-1}(I-P)t(Y) \qquad (2.3)$$

*is given by*

$$t(Y) = aE[\theta|Y] + (1-a)PE[\theta|Y]$$

*where*

$$a = a(Y) = \left[1 + \frac{H_1(Y)}{H_2(Y)}\right]^{1/2}$$

*and*

$$H_1(Y) = tr\left\{Var\left[\Phi^{-1/2}(I-P)\theta|Y\right]\right\}$$

$$H_2(Y) = E[\theta|Y]^T(I-P)^T\Phi^{-1}(I-P)E[\theta|Y].$$

Theorem 1 says that the projection of the estimates $t(Y)$ on $X$ matches the projection of the posterior mean of $\theta$ on $X$ and that the residual variance of the estimates about the regression surface matches the expected posterior residual variance of $\theta$ about the regression surface.

An important assumption of this result is that $\gamma$, $\sigma$, and $\phi$ are all known. Stern and Cressie (1999) give some direction on choosing the variance matrix $\Phi$. In particular, in order for $\Phi$ to be positive definite, $\phi \in (\phi_{min}, \phi_{max})$ where $\phi_{min} = \eta_1^{-1}$, $\phi_{max} = \eta_n^{-1}$, and $\eta_1 < 0 < \eta_n$, are, respectively, the smallest and largest eigenvalues of $M^{-1/2}CM^{1/2}$. When incorporating the conditional autoregressive model into a constrained Bayes context however, the results in Stern and Cressie (1999) assume $\phi$ is fixed.

## 2.3 Constrained Bayes with a General Covariance Matrix

We now consider the problem of fitting a more general small-area model using constrained Bayes methods. The basic model from Chapter 1 is

$$y_h = \theta_h + \epsilon_h, \ \{\epsilon_h\} \sim NID\left(0, \sigma_h^2\right)$$

$$\theta_h = x_h^T\gamma + \omega_h, \{\omega_h\} \sim NID\left(0, \ \tau_\omega^2\right),$$

with $\epsilon$ independent of $\omega$, which we extend to

$$y_h = \theta_h + \epsilon_h, \ \epsilon \sim N\left(0, V(\psi)\right)$$

$$\theta_h = x_h^T\gamma + \omega_h, \ \omega \sim N\left(0, \Sigma(\phi)\right). \tag{2.4}$$

In (2.4), $V(\psi)$ and $\Sigma(\phi)$ are covariance matrices dependent upon parameter vectors $\psi$ and $\phi$, respectively. We write $V = V(\psi)$ and $\Sigma = \Sigma(\phi)$ for compactness of notation.

Now, let $P = X\left(X^T\Sigma^{-1}X\right)^{-1}X^T\Sigma^{-1}$ denote the projection matrix that yields the predicted values for the generalized least squares regression on $X$ with error vector that has variance matrix $\Sigma$. We find the estimator $t(Y)$ that minimizes the posterior expected weighted squared error $E\left[(\theta - t(Y))^T\Sigma^{-1}(\theta - t(Y))|Y\right]$ subject to

$$E\left[P\theta|Y\right] = E\left[P|Y\right]t(Y) \tag{2.5}$$

and

$$E\left[\theta^T(I-P)^T\Sigma^{-1}(I-P)\theta|Y\right] = t(Y)^T E\left[(I-P)^T\Sigma^{-1}(I-P)|Y\right]t(Y). \tag{2.6}$$

Equation (2.5) is the analogue of (2.2), with $P$ allowed to depend on unknown parameters. Similarly, (2.6) is the analogue of (2.3).

### 2.3.1 General Constrained Bayes

We begin this section by stating the main result, which is an extension of Theorem 1 presented in Section 2.2.

**Result 1** *Suppose that* $Y|\theta \sim [\theta, V(\psi)]$ *and* $\theta \sim N[X\gamma, \Sigma(\phi)]$, *with* $V(\psi)$ *and* $\Sigma(\phi)$ *positive definite, and* $\gamma$, $\psi$, *and* $\phi$ *vectors of parameters on which the models depend. Writing* $\Sigma = \Sigma(\phi)$, *let* $P = X\left(X^T\Sigma^{-1}X\right)^{-1}X^T\Sigma^{-1}$ *denote the projection matrix that yields the predicted values for the generalized least squares regression on* $X$ *with error vector that has variance matrix* $\Sigma$, *and* $P^* = \left(X^T\Sigma^{-1}X\right)^{-1}X^T\Sigma^{-1}$. *Then, the* General Constrained Bayes (GCB) *estimates minimize the posterior expected weighted squared error*

$$E\left\{E\left[(\theta - t(Y))^T\Sigma^{-1}(\theta - t(Y))\,|\,Y,\phi\right]\,|\,Y\right\} \tag{2.7}$$

*subject to*

$$E\left[P^* E\left(\boldsymbol{\theta} \mid \boldsymbol{Y}, \phi\right) \mid \boldsymbol{Y}\right] = E\left[P^* \mid \boldsymbol{Y}\right] \boldsymbol{t}\left(\boldsymbol{Y}\right) \tag{2.8}$$

*and*

$$E\left\{E\left[\boldsymbol{\theta}^T\left(\boldsymbol{I} - \boldsymbol{P}\right)^T \boldsymbol{\Sigma}^{-1}\left(\boldsymbol{I} - \boldsymbol{P}\right)\boldsymbol{\theta} \mid \boldsymbol{Y}, \phi\right] \mid \boldsymbol{Y}\right\}$$
$$= \boldsymbol{t}\left(\boldsymbol{Y}\right)^T E\left[\left(\boldsymbol{I} - \boldsymbol{P}\right)^T \boldsymbol{\Sigma}^{-1}\left(\boldsymbol{I} - \boldsymbol{P}\right) \mid \boldsymbol{Y}\right] \boldsymbol{t}\left(\boldsymbol{Y}\right). \tag{2.9}$$

*The* GCB *estimates are given by* $\boldsymbol{t}\left(\hat{\lambda}_0, \hat{\boldsymbol{\lambda}}\right)$, *where*

$$\boldsymbol{t}\left(\lambda_0, \boldsymbol{\lambda}\right) = \left[2E\left(\boldsymbol{\Sigma}^{-1} \mid \boldsymbol{Y}\right) - 2\lambda_0 E\left(\left(\boldsymbol{I} - \boldsymbol{P}\right)^T \boldsymbol{\Sigma}^{-1}\left(\boldsymbol{I} - \boldsymbol{P}\right) \mid \boldsymbol{Y}\right)\right]^{-1} \times$$
$$\left[2E\left[\boldsymbol{\Sigma}^{-1} E\left(\boldsymbol{\theta} \mid \boldsymbol{Y}, \phi\right) \mid \boldsymbol{Y}\right] + \left[E\left(P^* \mid \boldsymbol{Y}\right)\right]^T \boldsymbol{\lambda}\right], \tag{2.10}$$

*and* $\left(\hat{\lambda}_0, \hat{\boldsymbol{\lambda}}\right)$ *are the solutions to*

$$0 = E\left\{E\left[\boldsymbol{\theta}^T\left(\boldsymbol{I} - \boldsymbol{P}\right)^T \boldsymbol{\Sigma}^{-1}\left(\boldsymbol{I} - \boldsymbol{P}\right)\boldsymbol{\theta} \mid \boldsymbol{Y}, \phi\right] \mid \boldsymbol{Y}\right\}$$
$$- E\left[\boldsymbol{t}\left(\lambda_0, \boldsymbol{\lambda}\right)\left(\boldsymbol{I} - \boldsymbol{P}\right)^T \boldsymbol{\Sigma}^{-1}\left(\boldsymbol{I} - \boldsymbol{P}\right)\boldsymbol{t}\left(\lambda_0, \boldsymbol{\lambda}\right) \mid \boldsymbol{Y}\right] \tag{2.11}$$

$$\boldsymbol{0} = E\left[PE\left(\boldsymbol{\theta} \mid \boldsymbol{Y}, \phi\right) \mid \boldsymbol{Y}\right] - E\left[P^* \boldsymbol{t}\left(\lambda_0, \boldsymbol{\lambda}\right) \mid \boldsymbol{Y}\right], \tag{2.12}$$

*respectively.*

Unlike the result of Theorem 1, the GCB estimator is not available analytically, and we resort to numerical methods in its computation. In particular, we use Lagrangian methods to minimize (2.7) subject to the constraints in (2.8) and (2.9). Writing $\boldsymbol{t} = \boldsymbol{t}\left(\boldsymbol{Y}\right)$, the Lagrangian is then

$$L = E\left\{E\left[\left(\boldsymbol{\theta} - \boldsymbol{t}\right)^T \boldsymbol{\Sigma}^{-1}\left(\boldsymbol{\theta} - \boldsymbol{t}\right) \mid \boldsymbol{Y}, \phi\right] \mid \boldsymbol{Y}\right\}$$
$$+ \lambda_0 E\left\{E\left[\boldsymbol{\theta}^T\left(\boldsymbol{I} - \boldsymbol{P}\right)^T \boldsymbol{\Sigma}^{-1}\left(\boldsymbol{I} - \boldsymbol{P}\right)\boldsymbol{\theta} \mid \boldsymbol{Y}, \phi\right] \mid \boldsymbol{Y}\right\}$$
$$- \lambda_0 \boldsymbol{t}^T E\left[\left(\boldsymbol{I} - \boldsymbol{P}\right)^T \boldsymbol{\Sigma}^{-1}\left(\boldsymbol{I} - \boldsymbol{P}\right) \mid \boldsymbol{Y}\right] \boldsymbol{t}$$
$$+ \boldsymbol{\lambda}^T \left\{E\left[P^* E\left(\boldsymbol{\theta} \mid \boldsymbol{Y}, \phi\right) \mid \boldsymbol{Y}\right] - E\left[P^* \mid \boldsymbol{Y}\right] \boldsymbol{t}\right\} \tag{2.13}$$

where $\boldsymbol{\lambda}$ is a $p \times 1$ vector.

Before continuing with the constrained minimization problem, we comment on the use of the constraint in (2.8) instead of

$$E\left[PE\left(\boldsymbol{\theta}\mid\boldsymbol{Y},\phi\right)\mid\boldsymbol{Y}\right]=E\left[\boldsymbol{P}\mid\boldsymbol{Y}\right]\boldsymbol{t}\left(\boldsymbol{Y}\right) \qquad (2.14)$$

Note that (2.14) represents a system of $m$ equations while (2.9) is a single equation. Thus, if we use the constraint in (2.14), we must solve for $m+1$ Lagrange multipliers. By using the constraint

$$E\left[\boldsymbol{P}^{*}E\left(\boldsymbol{\theta}\mid\boldsymbol{Y},\phi\right)\mid\boldsymbol{Y}\right]=E\left[\boldsymbol{P}^{*}\mid\boldsymbol{Y}\right]\boldsymbol{t}\left(\boldsymbol{Y}\right)$$

in place of (2.14), then we have reduced our problem to one in which we must find $p+1$ Lagrange multipliers, where $p$ is the number of covariates, typically much smaller than $m$.

To optimize $L$, take its partial derivatives with respect to $\boldsymbol{t}, \lambda_0$, and $\boldsymbol{\lambda}$ and set equal to zero:

$$0=\frac{\partial L}{\partial \boldsymbol{t}} = E\left[2\Sigma^{-1}\left(\boldsymbol{t}-E\left(\boldsymbol{\theta}\mid\boldsymbol{Y},\phi\right)\right)\mid\boldsymbol{Y}\right]$$
$$-2\lambda_0 E\left[\left(\boldsymbol{I}-\boldsymbol{P}\right)^{T}\Sigma^{-1}\left(\boldsymbol{I}-\boldsymbol{P}\right)\boldsymbol{t}\mid\boldsymbol{Y}\right]-\left[E\left(\boldsymbol{P}^{*}\mid\boldsymbol{Y}\right)\right]^{T}\boldsymbol{\lambda} \quad (2.15)$$

$$0=\frac{\partial L}{\partial \lambda_0} = E\left\{E\left[\boldsymbol{\theta}^{T}\left(\boldsymbol{I}-\boldsymbol{P}\right)^{T}\Sigma^{-1}\left(\boldsymbol{I}-\boldsymbol{P}\right)\boldsymbol{\theta}\mid\boldsymbol{Y},\phi\right]\mid\boldsymbol{Y}\right\}$$
$$-E\left[\boldsymbol{t}\left(\boldsymbol{I}-\boldsymbol{P}\right)^{T}\Sigma^{-1}\left(\boldsymbol{I}-\boldsymbol{P}\right)\boldsymbol{t}\mid\boldsymbol{Y}\right] \qquad (2.16)$$

$$0=\frac{\partial L}{\partial \boldsymbol{\lambda}} = E\left[\boldsymbol{P}^{*}E\left(\boldsymbol{\theta}\mid\boldsymbol{Y},\phi\right)\mid\boldsymbol{Y}\right]-E\left[\boldsymbol{P}^{*}\boldsymbol{t}\mid\boldsymbol{Y}\right] \qquad (2.17)$$

Solving (2.15) for $\boldsymbol{t}$ gives equation (2.10), which can then be substituted into (2.16) and (2.17). The resulting system of equations is (2.11) and (2.12), which in general requires a numerical solution for $\hat{\lambda}_0$ and $\hat{\boldsymbol{\lambda}}$.

Equations (2.10)–(2.12) require various posterior expectations, which in practice are approximated numerically. First note that, by standard computations (e.g., Gelman et al., 2004, p. 87),

$$\boldsymbol{\theta}\mid\boldsymbol{Y},\phi \sim N\left[\left(\Sigma^{-1}+\boldsymbol{V}^{-1}\right)^{-1}\left(\Sigma^{-1}\boldsymbol{X}\boldsymbol{\gamma}+\boldsymbol{V}^{-1}\boldsymbol{Y}\right),\left(\Sigma^{-1}+\boldsymbol{V}^{-1}\right)^{-1}\right]$$
$$= N\left(\boldsymbol{\mu},\boldsymbol{\Gamma}\right). \qquad (2.18)$$

Using the notation in (2.18), the posterior expectations in the expressions (2.11) and (2.12) can be simplified using

$$E\left(\boldsymbol{\theta} \mid \boldsymbol{Y}, \boldsymbol{\phi}\right) = \boldsymbol{\mu} \qquad (2.19)$$

and

$$
\begin{aligned}
E&\left[\boldsymbol{\theta}^T \left(\boldsymbol{I} - \boldsymbol{P}\right)^T \boldsymbol{\Sigma}^{-1} \left(\boldsymbol{I} - \boldsymbol{P}\right) \boldsymbol{\theta} \mid \boldsymbol{Y}, \boldsymbol{\phi}\right]\\
&= \; tr\left[\left(\boldsymbol{I} - \boldsymbol{P}\right)^T \boldsymbol{\Sigma}^{-1} \left(\boldsymbol{I} - \boldsymbol{P}\right) \boldsymbol{\Gamma}\right] + \boldsymbol{\mu}^T \left(\boldsymbol{I} - \boldsymbol{P}\right)^T \boldsymbol{\Sigma}^{-1} \left(\boldsymbol{I} - \boldsymbol{P}\right) \boldsymbol{\mu}, \quad (2.20)
\end{aligned}
$$

the latter due to Theorem 2.3, p. 49 in Hocking (1996).

The following is a numerical algorithm for computing the GCB using readily available software such as S-Plus or R:

1. Specify the model and obtain draws from the complete conditional distributions of all parameters in the model. For this we use WinBUGS (Lunn et al., 2000), which uses Gibbs sampling, possibly with a Metropolis step, to obtain realizations from each conditional distribution.

2. Import the draws into R (http://www.r-project.org/). Compute the following quantities:

   (a) $E\left(\boldsymbol{\Sigma}^{-1} \mid \boldsymbol{Y}\right)$. This is done by computing the function $\boldsymbol{\Sigma}^{-1}$ for each draw and taking the mean across all draws.

   (b) $E\left[\boldsymbol{\Sigma}^{-1}\left(\boldsymbol{I} - \boldsymbol{P}\right) \mid \boldsymbol{Y}\right]$. This is done by computing the function
   $$\boldsymbol{\Sigma}^{-1}\left[\boldsymbol{I} - \boldsymbol{X}\left(\boldsymbol{X}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{X}\right)^{-1} \boldsymbol{X}^T \boldsymbol{\Sigma}^{-1}\right]$$
   for each draw and taking the mean across all draws.

   (c) $E\left[\boldsymbol{\Sigma}^{-1} E\left(\boldsymbol{\theta} \mid \boldsymbol{Y}, \boldsymbol{\phi}\right) \mid \boldsymbol{Y}\right]$. This is done by computing $E\left(\boldsymbol{\theta} \mid \boldsymbol{Y}, \boldsymbol{\phi}\right) = \boldsymbol{\mu}$ from (2.19) for each draw from the Gibbs sampler. Then, $\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$ is computed for each iteration and the mean quantity is taken across all draws.

   (d) $E\left(\boldsymbol{P}^* \mid \boldsymbol{Y}\right) = E\left[\left(\boldsymbol{X}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{X}\right)^{-1} \boldsymbol{X}^T \boldsymbol{\Sigma}^{-1} \mid \boldsymbol{Y}\right]$. This is done by computing the function $\left(\boldsymbol{X}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{X}\right)^{-1} \boldsymbol{X}^T \boldsymbol{\Sigma}^{-1}$ for each draw and taking the mean across all draws.

(e) $E\left\{E\left[\boldsymbol{\theta}^{T}\left(\boldsymbol{I}-\boldsymbol{P}\right)^{T}\boldsymbol{\Sigma}^{-1}\left(\boldsymbol{I}-\boldsymbol{P}\right)\boldsymbol{\theta}\,|\,\boldsymbol{Y},\boldsymbol{\phi}\right]\,\Big|\,\boldsymbol{Y}\right\}$. The inner expectation can be computed for each draw from the Gibbs sampler by using the equality in (2.20). Taking the mean of this quantity over all draws gives an estimate of the outer expectation.

(f) $E\left[\left(\boldsymbol{I}-\boldsymbol{P}\right)^{T}\boldsymbol{\Sigma}^{-1}\left(\boldsymbol{I}-\boldsymbol{P}\right)|\boldsymbol{Y}\right]$. This is done by computing the function $\left(\boldsymbol{I}-\boldsymbol{P}\right)^{T}\boldsymbol{\Sigma}^{-1}\left(\boldsymbol{I}-\boldsymbol{P}\right)$, where

$$\boldsymbol{P}=\boldsymbol{X}\left(\boldsymbol{X}^{T}\boldsymbol{\Sigma}^{-1}\boldsymbol{X}\right)^{-1}\boldsymbol{X}^{T}\boldsymbol{\Sigma}^{-1},$$

for each draw and taking the mean over all draws.

(g) $E\left[\boldsymbol{P}^{*}E\left(\boldsymbol{\theta}\,|\,\boldsymbol{Y},\boldsymbol{\phi}\right)|\boldsymbol{Y}\right]$, computed as in (c) above.

3. Parts (a)–(d) of Step 2 allow computation of (2.10) up to the values of $\lambda_{0}$ and $\boldsymbol{\lambda}$. Parts (e)–(g) constitute the remaining terms in (2.11) and (2.12). Plugging in all of these values, solve for the Lagrangian multipliers $\lambda_{0}$ and $\boldsymbol{\lambda}$ (2.11) and (2.12), a system of $p+1$ equations. This nonlinear system may be solved in a number of ways. In particular, we have had success using the function `optim` in R, which uses a Nelder-Mead optimization algorithm, to minimize the sum of the squared estimation errors. With it, we compute (with some abuse of notation)

$$\left(\hat{\lambda}_{0},\hat{\boldsymbol{\lambda}}^{T}\right)^{T}=\arg\min_{\lambda_{0},\boldsymbol{\lambda}}\left[(2.11)^{2}+(2.12)^{T}(2.12)\right],$$

where (2.11) and (2.12) denote the right-hand sides of their respective equations.

4. Compute the GCB estimates by substituting the Lagrangian multipliers $\hat{\lambda}_{0}$ and $\hat{\boldsymbol{\lambda}}$ into (2.10) to obtain $\hat{\boldsymbol{t}}=\boldsymbol{t}\left(\hat{\lambda}_{0},\hat{\boldsymbol{\lambda}}\right)$.

## 2.4 Constrained Bayes Estimation on Different Scales

Recall that the CB estimators have the following properties:

$$\frac{1}{m}\sum_{h=1}^{m}\theta_{h}^{CB}=E\left[\bar{\boldsymbol{\theta}}|\boldsymbol{Y}\right] \tag{2.21}$$

and

$$\sum_{h=1}^{m} \left(\theta_h^{CB} - \overline{\theta^{CB}}\right)^2 = E\left[\sum_{h=1}^{n} \left(\theta_h - \bar{\theta}\right)^2 | \boldsymbol{Y}\right] \tag{2.22}$$

where $\boldsymbol{Y}$ is the data vector. It is of interest to know to what extent the properties of CB estimators are preserved under linear transformations. For example, we may wish to convert to a coarser set by aggregating small areas into larger regions. Do the CB properties hold for the coarser set? We show in this and the following section some *negative* results: CB properties are generally *not* preserved under linear transformations.

Let $\boldsymbol{\beta}^* = \boldsymbol{A}\boldsymbol{\theta}^{CB}$, where $\boldsymbol{A}$ is any full row rank $n \times m$ matrix, $n \leq m$. We begin by showing that $\{\beta_i^*\}_{i=1}^{n}$ generally does not satisfy the mean constraint in (2.21) nor the variance constraint in (2.22). First,

$$\frac{1}{n}\sum_{i=1}^{n} \beta_i^* = \frac{1}{n}\sum_{i=1}^{n} \left(\boldsymbol{A}\boldsymbol{\theta}^{CB}\right)_i$$

$$= \frac{1}{n}\boldsymbol{1}_n^T \boldsymbol{A}\boldsymbol{\theta}^{CB},$$

and

$$E\left[\bar{\beta}|\boldsymbol{Y}\right] = E\left[\frac{1}{n}\boldsymbol{1}_n^T \boldsymbol{A}\boldsymbol{\theta}|\boldsymbol{Y}\right]$$

$$= \frac{1}{n}\boldsymbol{1}_n^T \boldsymbol{A}E\left[\boldsymbol{\theta}|\boldsymbol{Y}\right].$$

Likewise,

$$\sum_{i=1}^{n} \left(\beta_i^* - \overline{\beta^*}\right)^2 = \left(\boldsymbol{\beta}^* - \overline{\beta^*}\boldsymbol{1}_n\right)^T \left(\boldsymbol{\beta}^* - \overline{\beta^*}\boldsymbol{1}_n\right)$$

$$= \left( A\theta^{CB} - \frac{1}{n} U_n A\theta^{CB} \right)^T \left( A\theta^{CB} - \frac{1}{n} U_n A\theta^{CB} \right)$$

$$= \left( A\theta^{CB} \right)^T A\theta^{CB} - \frac{1}{n} \left( A\theta^{CB} \right)^T U_n A\theta^{CB}$$

$$= \left( \theta^{CB} \right)^T \left( A^T A - \frac{1}{n} A^T U_n A \right) \theta^{CB}$$

where $U_n$ is an $n \times n$ matrix of 1's, and

$$E\left[ \sum_{i=1}^{n} \left( \beta_i - \bar{\beta} \right)^2 | \mathbf{Y} \right] = E\left[ \left( \boldsymbol{\beta} - \bar{\beta} \mathbf{1}_n \right)^T \left( \boldsymbol{\beta} - \bar{\beta} \mathbf{1}_n \right) | \mathbf{Y} \right]$$

$$= E\left[ \left( A\theta - \frac{1}{n} U_n A\theta \right)^T \left( A\theta - \frac{1}{n} U_n A\theta \right) | \mathbf{Y} \right]$$

$$= E\left[ \left( (A\theta)^T A\theta - \frac{1}{n} (A\theta)^T U_n A\theta \right) | \mathbf{Y} \right]$$

$$= E\left[ \theta^T \left( A^T A - \frac{1}{n} A^T U_n A \right) \theta | \mathbf{Y} \right]$$

So,

$$\frac{1}{n} \sum_{i=1}^{n} \beta_i^* = E\left[ \bar{\beta} | \mathbf{Y} \right]$$

$$\Leftrightarrow \frac{1}{n} \mathbf{1}_n^T A\theta^{CB} = \frac{1}{n} \mathbf{1}_n^T A E\left[ \theta | \mathbf{Y} \right], \tag{2.23}$$

where $\Leftrightarrow$ denotes if and only if. Furthermore,

$$\sum_{i=1}^{n} \left( \beta_i^* - \overline{\beta^*} \right)^2 = E\left[ \sum_{i=1}^{n} \left( \beta_i - \bar{\beta} \right)^2 | \mathbf{Y} \right]$$

$$\Leftrightarrow \left( \theta^{CB} \right)^T \left( A^T A - \frac{1}{n} A^T U_n A \right) \theta^{CB} = E\left[ \theta^T \left( A^T A - \frac{1}{n} A^T U_n A \right) \theta | \mathbf{Y} \right]$$

$$\Leftrightarrow tr\left\{ \left( \theta^{CB} \right)^T \left( A^T A - \frac{1}{n} A^T U_n A \right) \theta^{CB} \right\} = tr\left\{ E\left[ \theta^T \left( A^T A - \frac{1}{n} A^T U_n A \right) \theta | \mathbf{Y} \right] \right\}$$

$$\Leftrightarrow tr\left\{ \left( A^T A - \frac{1}{n} A^T U_n A \right) \theta^{CB} \left( \theta^{CB} \right)^T \right\} = tr\left\{ E\left[ \left( A^T A - \frac{1}{n} A^T U_n A \right) \theta\theta^T | \mathbf{Y} \right] \right\}$$

$$\Leftrightarrow tr\left\{ \left( A^T A - \frac{1}{n} A^T U_n A \right) \theta^{CB} \left( \theta^{CB} \right)^T \right\} = tr\left\{ \left( A^T A - \frac{1}{n} A^T U_n A \right) E\left( \theta\theta^T | \mathbf{Y} \right) \right\}$$

$$\Leftrightarrow tr\left\{ A^* \theta^{CB} \left( \theta^{CB} \right)^T \right\} = tr\left\{ A^* E\left( \theta\theta^T | \mathbf{Y} \right) \right\}. \tag{2.24}$$

The conditions in (2.23) and (2.24) will hold in some cases, but in general, they do

not as illustrated in the following numerical counterexample.

Let $y|\theta \sim N(\theta, \sigma^2 I)$ and $\theta \sim N(\mu, \tau^2 I)$. Then, straightforward Bayes and

Constrained Bayes manipulations will show that

$$E[\theta|y] = \frac{\sigma^2 \mu + \tau^2 y}{\sigma^2 + \tau^2},$$

$$E[\theta\theta^T|y] = \left(\frac{\tau^2 \sigma^2}{\sigma^2 + \tau^2}\right) I_m + E[\theta|y] E[\theta|y]^T,$$

$$\theta^{CB} = \left(a I_m + \frac{1-a}{m} U_m\right) E[\theta|y], \text{ and}$$

$$a = \left(1 + \frac{(m-1)\frac{\tau^2 \sigma^2}{\sigma^2 + \tau^2}}{E[\theta|y]^T \left(I_m - \frac{1}{m} U_m\right) E[\theta|y]}\right)^{1/2}.$$

Suppose $\mu = (10\ 20\ 30\ 40)^T$, $\sigma^2 = 5$, $\tau^2 = 10$, and $y = (18\ 16\ 27\ 34)^T$. Then,

$$E[\theta|y] = \left(15\frac{1}{3}\ 17\frac{1}{3}\ 28\ 36\right)^T,$$

$$\theta^{CB} = \left(a I_4 + \frac{1-a}{4} U_4\right)\left(15\frac{1}{3}\ 17\frac{1}{3}\ 28\ 36\right)^T,$$

$$E[\theta\theta^T|y] = \begin{bmatrix} 238\frac{4}{9} & 265\frac{7}{9} & 429\frac{1}{3} & 552 \\ 265\frac{7}{9} & 303\frac{7}{9} & 485\frac{1}{3} & 624 \\ 429\frac{1}{3} & 485\frac{1}{3} & 787\frac{1}{3} & 1008 \\ 552 & 624 & 1008 & 1299\frac{1}{3} \end{bmatrix}, \text{ and}$$

$$\theta^{CB}(\theta^{CB})^T = \left(a I_4 + \frac{1-a}{4} U_4\right) \begin{bmatrix} 235\frac{1}{9} & 265\frac{7}{9} & 429\frac{1}{3} & 552 \\ 265\frac{7}{9} & 300\frac{4}{9} & 485\frac{1}{3} & 624 \\ 429\frac{1}{3} & 485\frac{1}{3} & 784 & 1008 \\ 552 & 624 & 1008 & 1296 \end{bmatrix} \left(a I_4 + \frac{1-a}{4} U_4\right)^T,$$

where

$$a = \left(\frac{521}{503}\right)^{1/2}.$$

If

$$A = \begin{bmatrix} 10 & 1 & 0 & 0 \\ 0 & 0 & 1 & .1 \end{bmatrix},$$

then

$$\frac{1}{m}\mathbf{1}_m^T A\boldsymbol{\theta}^{CB} = 100\frac{3339}{10000} \neq 101\frac{2}{15} = \frac{1}{m}\mathbf{1}_m^T AE\left[\boldsymbol{\theta}|\boldsymbol{y}\right], \text{ and}$$

$$tr\left\{A^*\boldsymbol{\theta}^{CB}\left(\boldsymbol{\theta}^{CB}\right)^T\right\} = 19491\frac{3}{20} \neq 20152\frac{3}{4} = tr\left\{A^*E\left(\boldsymbol{\theta}\boldsymbol{\theta}^T|\boldsymbol{y}\right)\right\}.$$

In general then, $\boldsymbol{\beta}^{CB} \neq A\boldsymbol{\theta}^{CB}$. In the next section, we explore further relationships between $\boldsymbol{\beta}^{CB}$ and $A\boldsymbol{\theta}^{CB}$.

### 2.4.1 Search for alternative criteria

In this section, we will assume that we are estimating $\boldsymbol{\theta}$. In addition to finding a good ensemble estimate $\tilde{\boldsymbol{\theta}}$ of $\boldsymbol{\theta}$, we would like a good ensemble estimate $\tilde{\boldsymbol{\beta}} = A\tilde{\boldsymbol{\theta}}$ for $\boldsymbol{\beta} = A\boldsymbol{\theta}$. We already know that $\tilde{\boldsymbol{\theta}} = \boldsymbol{\theta}^{CB}$ has good properties with respect to mean and variability. But as shown in the previous section, $\boldsymbol{\beta}^* = A\boldsymbol{\theta}^{CB}$ does not necessarily share these properties. This section focuses on the search for an ensemble estimate $\tilde{\boldsymbol{\theta}}$ such that $\tilde{\boldsymbol{\beta}} = A\tilde{\boldsymbol{\theta}}$ satisfies, approximately, the mean and variance constraints in (2.21) and (2.22).

We begin by showing that if the variance constraint in (2.22) is satisfied, then the mean constraint (2.21) is automatically satisfied. Using a Lagrange multiplier argument to perform constrained optimization, we minimize the mean squared error

$$E\left[\sum_{h=1}^{m}\left(\theta_h - \tilde{\theta}_h\right)^2 |\boldsymbol{Y}\right]$$

subject to

$$\sum_{h=1}^{m} \left(\tilde{\theta}_h - \bar{\tilde{\theta}}\right)^2 = E\left[\sum_{h=1}^{m} \left(\theta_h - \bar{\theta}\right)^2 | \mathbf{Y}\right]. \qquad (2.25)$$

The Lagrangian is then

$$
\begin{aligned}
L =\ & E\left[\left(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}\right)^T \left(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}\right) | \mathbf{Y}\right] \\
& + \lambda\left\{\left(\tilde{\boldsymbol{\theta}} - \frac{1}{m}U_m\tilde{\boldsymbol{\theta}}\right)^T \left(\tilde{\boldsymbol{\theta}} - \frac{1}{m}U_m\tilde{\boldsymbol{\theta}}\right) - E\left[\left(\boldsymbol{\theta} - \frac{1}{m}U_m\boldsymbol{\theta}\right)^T \left(\boldsymbol{\theta} - \frac{1}{m}U_m\boldsymbol{\theta}\right) | \mathbf{Y}\right]\right\}
\end{aligned}
$$

$$(2.26)$$

To minimize, we take the partial derivative with respect to $\tilde{\boldsymbol{\theta}}$ and set equal to zero:

$$
\begin{aligned}
\frac{\partial L}{\partial \tilde{\boldsymbol{\theta}}} = 2\left(\tilde{\boldsymbol{\theta}} - E\left[\boldsymbol{\theta}|\mathbf{Y}\right]\right) + 2\lambda\left(I_m - \frac{1}{m}U_m\right)\tilde{\boldsymbol{\theta}} &= 0 \\
\Leftrightarrow \frac{1}{m}U_m\left(\tilde{\boldsymbol{\theta}} - E\left[\boldsymbol{\theta}|\mathbf{Y}\right]\right) + \lambda\frac{1}{m}U_m\left(I_m - \frac{1}{m}U_m\right)\tilde{\boldsymbol{\theta}} &= 0 \\
\Leftrightarrow \frac{1}{m}U_m\left(\tilde{\boldsymbol{\theta}} - E\left[\boldsymbol{\theta}|\mathbf{Y}\right]\right) + \lambda\frac{1}{m}U_m\tilde{\boldsymbol{\theta}} - \lambda\frac{1}{m}U_m\tilde{\boldsymbol{\theta}} &= 0 \\
\Leftrightarrow \frac{1}{m}U_m\left(\tilde{\boldsymbol{\theta}} - E\left[\boldsymbol{\theta}|\mathbf{Y}\right]\right) &= 0 \\
\Leftrightarrow \frac{1}{m}\sum_{i=1}^{m}\tilde{\theta}_i &= E\left[\bar{\boldsymbol{\theta}}|\mathbf{Y}\right].
\end{aligned}
$$

Thus, we focus our attention on the variance constraint.

Equation (2.24) gives us some guidance as to additional constraints we might consider in order to find an estimator $\tilde{\boldsymbol{\theta}}$ that will ensure

$$\sum_{i=1}^{n} \left(\tilde{\beta}_i - \bar{\tilde{\beta}}\right)^2 = E\left[\sum_{i=1}^{n} \left(\beta_i - \bar{\beta}\right)^2 | \mathbf{Y}\right]. \qquad (2.27)$$

For example, it is clear from (2.24) that if

$$\tilde{\boldsymbol{\theta}}\tilde{\boldsymbol{\theta}}^T = E\left(\boldsymbol{\theta}\boldsymbol{\theta}^T|\mathbf{Y}\right), \qquad (2.28)$$

then (2.27) will hold. Examining this condition, we notice a difficulty, namely that every element of the two matrices must be equal. Consider again a simple example where $Y|\theta \sim N(\theta, \sigma^2 I)$ and $\theta \sim N(\mu, \tau^2 I)$. Let $\theta = (\theta_1, \theta_2)$. Then,

$$E\left[\theta\theta^T|Y\right] = E\left[\begin{pmatrix} \theta_1^2 & \theta_1\theta_2 \\ \theta_1\theta_2 & \theta_2^2 \end{pmatrix}|Y\right].$$

So, we have to match three quantities:

$$\tilde{\theta}_1^2 \equiv E\left[\theta_1^2|Y\right],$$

$$\tilde{\theta}_2^2 \equiv E\left[\theta_2^2|Y\right], \text{ and}$$

$$\tilde{\theta}_1\tilde{\theta}_2 \equiv E\left[\theta_1\theta_2|Y\right]. \tag{2.29}$$

But, as long as $\sigma^2, \tau^2 > 0$,

$$E\left[\theta_1^2|Y\right] = \frac{\tau^2\sigma^2}{\sigma^2 + \tau^2} + \left(\frac{\sigma^2\mu_1 + \tau^2 y_1}{\sigma^2 + \tau^2}\right)^2,$$

$$E\left[\theta_2^2|Y\right] = \frac{\tau^2\sigma^2}{\sigma^2 + \tau^2} + \left(\frac{\sigma^2\mu_2 + \tau^2 y_2}{\sigma^2 + \tau^2}\right)^2, \text{ and}$$

$$E\left[\theta_1\theta_2|Y\right] = \left(\frac{\sigma^2\mu_1 + \tau^2 y_1}{\sigma^2 + \tau^2}\right)\left(\frac{\sigma^2\mu_2 + \tau^2 y_2}{\sigma^2 + \tau^2}\right)$$

$$\neq \sqrt{E\left[\theta_1^2|Y\right]E\left[\theta_2^2|Y\right]},$$

as would be required by the assignments in (2.29). More generally, the constraint in (2.28) would always require $E\left[\theta_i\theta_j|Y\right] = \sqrt{E\left[\theta_i^2|Y\right]E\left[\theta_j^2|Y\right]}$. But, using Cauchy-Schwartz,

$$E\left[\theta_i\theta_j|Y\right] \leq |E\left[\theta_i\theta_j|Y\right]| \leq \sqrt{E\left[\theta_i^2|Y\right]E\left[\theta_j^2|Y\right]},$$

with $|E\left[\theta_i\theta_j|Y\right]| = \sqrt{E\left[\theta_i^2|Y\right]E\left[\theta_j^2|Y\right]}$ if and only if one of $\theta_i|Y, \theta_j|Y = 0$ or $\theta_j|Y = c\theta_i|Y$ for some constant $c$. Thus, in most useful situations, the constraint in (2.28) will not be attainable.

There are two directions this investigation might take from this point: one, we can try to find a set of estimators $\tilde{\theta}$ of $\theta$ for which the constraints hold *approximately*, or two, we can try to find a class of linear transformations represented by $A$ for which the constraints hold. In the last chapter of this dissertation, we discuss some candidates for $\tilde{\theta}$ and $A$ that have been investigated and shown to be inadequate and some candidates that will be investigated in future work.

## 2.5 Bias of $F_{\theta^B}(k)$ and $F_{\theta^{CB}}(k)$

As demonstrated earlier, $F_{\theta^{CB}}(k)$, the empirical distribution of the constrained Bayes estimators of $\theta$, appears to be a better estimator of $F_\theta^B(k)$, the posterior mean of the empirical distribution of $\theta$ than is $F_{\theta^B}(k)$, the empirical distribution of the Bayes estimators of $\theta$. We seek an analytical demonstration of this observation under a normal model and use the results as a point of comparison for the relationship among $F_\beta^B(k)$, $F_{\beta^*}(k)$, $F_{\beta'}(k)$, where $\beta^*$ is defined in Section 2.4 and $\beta' = A\theta^B$. That is, we have demonstrated that $\beta^* = A\theta^{CB}$ is *not* CB, but does it at least have better ensemble estimation properties than $\beta'$?

Assume the conditional model

$$
\begin{aligned}
y|\theta &\sim N\left(\theta, \sigma^2 I\right) \\
\theta &\sim N\left(\mu, \tau^2 I\right),
\end{aligned}
\tag{2.30}
$$

which leads to a joint model of

$$
\begin{pmatrix} y \\ \theta \end{pmatrix} \sim N\left[ \begin{pmatrix} \mu \\ \mu \end{pmatrix}, \begin{pmatrix} \left(\tau^2 + \sigma^2\right) I & \tau^2 I \\ \tau^2 I & \tau^2 I \end{pmatrix} \right].
\tag{2.31}
$$

Define the bias of $F_{\theta B}(k)$ as

$$
\begin{aligned}
E\left[F_{\theta B}(k) - F_\theta^B(k)\right] &= E\left\{\frac{1}{m}\sum_{h=1}^m I\left\{E\left(\theta_h|\boldsymbol{y}\right) \le k\right\} - E\left[\frac{1}{m}\sum_{h=1}^m I\left\{\theta_h \le k\right\}|\boldsymbol{y}\right]\right\} \\
&= \frac{1}{m}\sum_{h=1}^m \left\{E\left[I\left\{E\left(\theta_h|\boldsymbol{y}\right)\right\} \le k\right] - E\left[E\left[I\left\{\theta_h \le k\right\}|\boldsymbol{y}\right]\right]\right\} \\
&= \frac{1}{m}\sum_{h=1}^m \left\{P\left[E\left(\theta_h|\boldsymbol{y}\right) \le k\right] - E\left[I\left\{\theta_h \le k\right\}\right]\right\} \\
&= \frac{1}{m}\sum_{h=1}^m \left\{P\left[\frac{\tau^2}{\sigma^2+\tau^2}y_h + \frac{\sigma^2}{\sigma^2+\tau^2}\mu_h \le k\right] - E\left[I\left\{\theta_h \le k\right\}\right]\right\} \\
&= \frac{1}{m}\sum_{h=1}^m \left\{P\left[y_h \le \frac{k - \frac{\sigma^2}{\sigma^2+\tau^2}\mu_h}{\frac{\tau^2}{\sigma^2+\tau^2}}\right] - E\left[I\left\{\theta_h \le k\right\}\right]\right\} \\
&= \frac{1}{m}\sum_{h=1}^m \left\{P\left[z_h \le \frac{\frac{k - \frac{\sigma^2}{\sigma^2+\tau^2}\mu_h}{\frac{\tau^2}{\sigma^2+\tau^2}} - \mu_h}{\sqrt{\sigma^2+\tau^2}}\right] - E\left[I\left\{\theta_h \le k\right\}\right]\right\} \\
&= \frac{1}{m}\sum_{h=1}^m \left\{P\left[z_h \le \frac{k - \mu_h}{\tau\sqrt{\frac{\tau^2}{\sigma^2+\tau^2}}}\right] - P\left[z_h \le \frac{k - \mu_h}{\tau}\right]\right\}, \quad (2.32)
\end{aligned}
$$

where $\{z_h\}_{h=1}^m$ are NID$(0,1)$.

For the case of (2.32) in which $\mu_h \equiv \mu$ for all $h$, we get

$$
\begin{aligned}
(2.32) &= P\left[z \le \frac{k - \mu}{\tau\sqrt{\frac{\tau^2}{\sigma^2+\tau^2}}}\right] - P\left[z \le \frac{k - \mu}{\tau}\right] \quad (2.33) \\
&\ne 0,
\end{aligned}
$$

demonstrating that, for this special case $F_{\theta B}$ is not even asymptotically unbiased for $F_\theta^B$. An illustration of this is shown in Figure 2.1 for which we let $\boldsymbol{\mu} = \boldsymbol{0}$. The nine different graphs show the bias for different ratios of $\sigma^2/\tau^2$, with $\tau^2$ fixed at 1. The bias is much higher as $\sigma^2$ gets larger. This is because the scale factor $\sqrt{\tau^2/(\sigma^2+\tau^2)}$ get smaller, shrinking the first probability in (2.32) closer to zero.

Figure 2.1: Bias of $F_{\theta^B}(k)$. The value of $\tau^2$ is fixed at 1, $\sigma^2 = 0.05, 0.1, 0.2, 0.5, 1, 2, 5, 10, 20$, and $m = 1000$. Here, we assume $\boldsymbol{\mu}$ is the 0 vector.

The bias of $F_{\theta^{CB}}(k)$ is more complex because the form of $\{\boldsymbol{\theta}^{CB}\}$ contains several functions of $\{\boldsymbol{y}\}$, the probabilistic component. With $a(\boldsymbol{y})$ defined as in Chapter 1,

$$E\left[F_{\theta^{CB}}(k)\right] = E\left\{\frac{1}{m}\sum_{h=1}^{m} I\left\{\theta_h^{CB} \leq k\right\}\right\}$$

$$= E\left\{\frac{1}{m}\sum_{h=1}^{m} I\left\{\bar{\theta}^B + a\left(\mathbf{y}\right)\left(\theta_h^B - \bar{\theta}^B\right) \le k\right\}\right\}. \qquad (2.34)$$

In (2.34), $\bar{\theta}^B$, $a\left(\mathbf{y}\right)$, and $\theta_h^B$ are all functions of $\{\mathbf{y}\}$. To begin, we find an expression

for $a\left(\mathbf{y}\right)$ under the model (2.31):

$$
\begin{aligned}
a\left(\mathbf{y}\right) &= \left(1 + \frac{tr\left\{Var\left(\boldsymbol{\theta} - \bar{\theta}\mathbf{1}\mid\mathbf{Y}\right)\right\}}{\sum_{h=1}^{m}\left(\theta_h^B - \bar{\theta}^B\right)^2}\right)^{1/2} \\
&= \left(1 + \frac{tr\left\{Var\left[\left(\mathbf{I} - \frac{1}{n}U\right)\boldsymbol{\theta}\mid\mathbf{Y}\right]\right\}}{\sum_{h=1}^{m}\left[E\left(\theta_h\mid\mathbf{Y}\right) - \frac{1}{n}\sum_{h=1}^{m}E\left(\theta_h\mid\mathbf{Y}\right)\right]^2}\right)^{1/2} \\
&= \left(1 + \frac{tr\left\{\left(\mathbf{I} - \frac{1}{n}U\right)\frac{\tau^2\sigma^2}{\tau^2+\sigma^2}\left(\mathbf{I} - \frac{1}{n}U\right)\right\}}{\sum_{h=1}^{m}\left[\frac{1}{\tau^2+\sigma^2}\left(\tau^2 y_h + \sigma^2\mu_h\right) - \frac{1}{n}\sum_{h=1}^{m}\frac{1}{\tau^2+\sigma^2}\left(\tau^2\bar{y} + \sigma^2\bar{\mu}\right)\right]^2}\right)^{1/2} \\
&= \left(1 + \frac{\frac{\tau^2\sigma^2}{\tau^2+\sigma^2}\left(\mathbf{I} - \frac{1}{n}U\right)}{\left(\frac{1}{\tau^2+\sigma^2}\right)^2\sum_{h=1}^{m}\left[\tau^2\left(y_h - \bar{y}\right) + \sigma^2\left(\mu_h - \bar{\mu}\right)\right]^2}\right)^{1/2} \\
&= \left(1 + \frac{\frac{\tau^2\sigma^2}{\tau^2+\sigma^2}\left(n - 1\right)}{\left(\frac{1}{\tau^2+\sigma^2}\right)^2\sum_{h=1}^{m}\left[\tau^2\left(y_h - \bar{y}\right) + \sigma^2\left(\mu_h - \bar{\mu}\right)\right]^2}\right)^{1/2} \\
&= \left(1 + \frac{\frac{\tau^2\sigma^2}{\tau^2+\sigma^2}}{\left(\frac{1}{\tau^2+\sigma^2}\right)^2 D}\right)^{1/2} \\
&= \left(1 + \frac{\tau^2\sigma^2}{\frac{1}{\tau^2+\sigma^2}D}\right)^{1/2}, \qquad (2.35)
\end{aligned}
$$

where

$$D = \tau^4\sum_{h=1}^{m}\frac{\left(y_h - \bar{y}\right)^2}{n - 1} + \sigma^4\sum_{h=1}^{m}\frac{\left(\mu_h - \bar{\mu}\right)^2}{n - 1} + 2\tau^2\sigma^2\sum_{h=1}^{m}\frac{\left(y_h - \bar{y}\right)\left(\mu_h - \bar{\mu}\right)}{n - 1}.$$

The first term of $D$ is an estimate of $\tau^4\left(\tau^2 + \sigma^2\right)$. Likewise, the second and third

terms of $D$ are representations of the variability in $\{\boldsymbol{\mu}\}$ and the covariability of $\{\mathbf{y}\}$

and $\{\boldsymbol{\mu}\}$, respectively. To simplify the problem of understanding the bias of $F_{\theta^{CB}}\left(k\right)$,

we assume $\{\boldsymbol{\mu}\}$ is a constant vector, and replace $a\left(\mathbf{y}\right)$ with its limit as $m \to \infty$. Then,

the second and third terms of $D$ become 0 and (2.35), in its limit, becomes

$$a_{\lim} = \left(1 + \frac{\sigma^2}{\tau^2}\right)^{1/2}. \qquad (2.36)$$

Using this expression for $a(\boldsymbol{y})$, we have

$$E\left[F_{\theta^{CB}}(k) - F_{\theta}^{B}(k)\right] = \frac{1}{m}\sum_{h=1}^{m}\left\{P\left[\bar{\theta}^B + a_{\lim}\left(\theta_h^B - \bar{\theta}^B\right) \leq k\right] - P\left[z_h \leq \frac{k - \mu_h}{\tau}\right]\right\}. \qquad (2.37)$$

Now, the first probability in (2.37) is

$$P\left[\bar{\theta}^B + a_{\lim}\left(\theta_h^B - \bar{\theta}^B\right) \leq k\right] = P\left[\frac{\tau^2}{\sigma^2+\tau^2}\sum_{h=1}^{m}y_h + \frac{\sigma^2}{\sigma^2+\tau^2}\mu\right.$$
$$+ a_{\lim}\frac{\tau^2}{\sigma^2+\tau^2}y_h + a_{\lim}\frac{\sigma^2}{\sigma^2+\tau^2}\mu$$
$$\left. - \frac{a_{\lim}}{m}\frac{\tau^2}{\sigma^2+\tau^2}\sum_{h=1}^{m}y_h - a_{\lim}\frac{\sigma^2}{\sigma^2+\tau^2}\mu \leq k\right]$$
$$= P\left[\frac{1}{m}\sum_{h=1}^{m}y_h + a_{\lim}y_h - \frac{a_{\lim}}{m}\sum_{h=1}^{m}y_h \leq \frac{k - \frac{\sigma^2}{\sigma^2+\tau^2}\mu}{\frac{\tau^2}{\sigma^2+\tau^2}}\right]. \qquad (2.38)$$

Furthermore,

$$\frac{1}{m}\sum_{h=1}^{m}y_h + a_{\lim}y_h - \frac{a_{\lim}}{m}\sum_{h=1}^{m}y_h \sim N\left[\mu, (\sigma^2+\tau^2)\frac{1 - a_{\lim}^2 + ma_{\lim}^2}{m}\right].$$

So,

$$P\left[\bar{\theta}^B + a_{\lim}\left(\theta_h^B - \bar{\theta}^B\right) \leq k\right] = P\left[z_h \leq \frac{\frac{k - \frac{\sigma^2}{\sigma^2+\tau^2}\mu}{\frac{\tau^2}{\sigma^2+\tau^2}} - \mu}{\sqrt{(\sigma^2+\tau^2)\frac{1 - a_{\lim}^2 + ma_{\lim}^2}{m}}}\right]$$
$$= P\left[z_h \leq \frac{k - \mu}{\frac{\tau^2}{\sqrt{\sigma^2+\tau^2}}}\sqrt{\frac{m}{1 - a_{\lim}^2 + ma_{\lim}^2}}\right], \quad (2.39)$$

and using the expression for $a_{\lim}$ in (2.36), we get

$$
\begin{aligned}
E\left[F_{\theta CB}\left(k\right)-F_{\theta}^{B}\left(k\right)\right] &= \frac{1}{m}\sum_{h=1}^{m}\left\{P\left[z_{h}\leq\frac{k-\mu}{\tau\sqrt{1-\frac{\sigma^{2}}{m(\sigma^{2}+\tau^{2})}}}\right]-P\left[z_{h}\leq\frac{k-\mu}{\tau}\right]\right\} \\
&= P\left[z\leq\frac{k-\mu}{\tau\sqrt{1-\frac{\sigma^{2}}{m(\sigma^{2}+\tau^{2})}}}\right]-P\left[z\leq\frac{k-\mu}{\tau}\right] \quad (2.40)\\
&\rightarrow \quad 0 \text{ as } m\rightarrow\infty.
\end{aligned}
$$

Equation (2.40) indicates that the bias of $F_{\theta CB}$ is relatively high when $\sigma^{2}$ is large relative to $\tau^{2}$. However, as $m$ increases, the bias gets smaller. In fact, (2.40) shows that $F_{\theta CB}$ is asymptotically unbiased for $F_{\theta}^{B}$ in the special case where $\mu_{h}\equiv\mu$ in contrast to (2.33) which shows that $F_{\theta B}$ is always biased. Figures 2.1 and 2.2 show the bias of $F_{\theta B}$ and $F_{\theta CB}$, respectively, for different values of the ratio $\sigma^{2}/\tau^{2}$ and assuming $\boldsymbol{\mu}=\mathbf{0}$. From the figures we see that when the variance $(\tau^{2})$ of the prior distribution of $\boldsymbol{\theta}$ is large compared to the variance of the measurement error $(\sigma^{2})$, $\sqrt{1-\frac{\sigma^{2}}{m(\sigma^{2}+\tau^{2})}}$ is close to 1 and there is very little bias for all values of $k$ as shown in the upper left panel. Note that the bias of $F_{\theta CB}$ is smaller than the bias of $F_{\theta B}$ by three orders of magnitude.

We turn now to the bias of $F_{\beta'}$ and $F_{\beta^{*}}$. First, using $A_{ih}$ to denote the element of $\boldsymbol{A}$ in the $i$th row and $h$th column,

$$
\begin{aligned}
E\left[F_{\beta'}\left(k\right)-F_{\beta}^{B}\left(k\right)\right] &= E\left\{\frac{1}{n}\sum_{i=1}^{n}I\left\{\beta_{i}'\leq k\right\}-E\left[\frac{1}{n}\sum_{i=1}^{n}I\left\{\beta_{i}\leq k\right\}|\boldsymbol{y}\right]\right\} \\
&= E\left\{\frac{1}{n}\sum_{i=1}^{n}I\left\{\sum_{h=1}^{m}A_{ih}E\left(\theta_{h}|\boldsymbol{y}\right)\leq k\right\}\right. \\
&\quad \left.-E\left[\frac{1}{n}\sum_{i=1}^{n}I\left\{\sum_{h=1}^{m}A_{ih}\theta_{h}\leq k\right\}|\boldsymbol{y}\right]\right\}
\end{aligned}
$$

Figure 2.2: Bias of $F_{\theta CB}(k)$. The value of $\tau^2$ is fixed at 1, $\sigma^2 = 0.05, 0.1, 0.2, 0.5, 1, 2, 5, 10, 20$, and $m = 1000$.

$$= \frac{1}{n}\sum_{i=1}^{n}\left\{ P\left[\sum_{h=1}^{m}A_{ih}E\left(\theta_h|y\right) \le k\right] - P\left[\sum_{h=1}^{m}A_{ih}\theta_h \le k\right]\right\}$$

$$= \frac{1}{n}\sum_{i=1}^{n}\left\{ P\left[\frac{\tau^2}{\sigma^2+\tau^2}\sum_{h=1}^{m}A_{ih}y_h + \frac{\sigma^2}{\sigma^2+\tau^2}\sum_{h=1}^{m}A_{ih}\mu_h \le k\right]\right.$$

$$\left. -P\left[\sum_{h=1}^{m}A_{ih}\theta_h \le k\right]\right\}$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left\{ P \left[ \sum_{h=1}^{m} A_{ih} y_h \leq \frac{k - \frac{\sigma^2}{\sigma^2 + \tau^2} \sum_{h=1}^{m} A_{ih} \mu_h}{\frac{\tau^2}{\sigma^2 + \tau^2}} \right] \right.$$

$$\left. - P \left[ \sum_{h=1}^{m} A_{ih} \theta_h \leq k \right] \right\}$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left\{ P \left[ z_i \leq \frac{k - \sum_{h=1}^{m} A_{ih} \mu_h}{\frac{\tau^2}{\sqrt{(\sigma^2 + \tau^2)}} \sqrt{\sum_{h=1}^{m} A_{ih}^2}} \right] \right.$$

$$\left. - P \left[ z_i \leq \frac{k - \sum_{h=1}^{m} A_{ih} \mu_h}{\tau \sqrt{\sum_{h=1}^{m} A_{ih}^2}} \right] \right\}. \qquad (2.41)$$

Equation (2.41) indicates that the bias of $F_{\beta'}$ is relatively high when $\sigma^2$ is large relative to $\tau^2$.

In deriving the bias of $F_{\beta^*}$, we make the same simplifications as we did when deriving the bias of $F_{\theta^{CB}}$. By following steps similar to those shown in deriving (2.38), (2.39), and (2.41), we find that

$$E \left[ F_{\beta^*}(k) - F_{\beta}^B(k) \right] = E \left\{ \frac{1}{n} \sum_{i=1}^{n} I \left\{ \beta_i^* \leq k \right\} - E \left[ \frac{1}{n} \sum_{i=1}^{n} I \left\{ \beta_i \leq k \right\} | \boldsymbol{y} \right] \right\}$$

$$= E \left\{ \frac{1}{n} \sum_{i=1}^{n} I \left\{ \sum_{h=1}^{m} A_{ih} \theta_h^{CB} \leq k \right\} \right.$$

$$\left. - E \left[ \frac{1}{n} \sum_{i=1}^{n} I \left\{ \sum_{h=1}^{m} A_{ih} \theta_h \leq k \right\} | \boldsymbol{y} \right] \right\}$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left\{ P \left[ \sum_{h=1}^{m} A_{ih} \theta_h^{CB} \leq k \right] \right.$$

$$\left. - P \left[ \sum_{h=1}^{m} A_{ih} \theta_h \leq k \right] \right\} \qquad (2.42)$$

Now, the first probability in (2.42) is

$$P \left[ \sum_{h=1}^{m} A_{ih} \theta_h^{CB} \leq k \right] = P \left[ \sum_{h=1}^{m} A_{ih} \left( \bar{\theta}^B + a_{\lim} \left( \theta_h^B - \bar{\theta}^B \right) \right) \leq k \right]$$

$$= P \left[ A_{i \cdot} \frac{\tau^2}{\sigma^2 + \tau^2} \frac{1}{m} \sum_{h=1}^{m} A_{ih} y_h + A_{i \cdot} \frac{\sigma^2}{\sigma^2 + \tau^2} \mu \right.$$

$$+ a_{\lim} \frac{\tau^2}{\sigma^2 + \tau^2} \sum_{h=1}^{m} A_{ih} y_h + a_{\lim} A_{i\cdot} \frac{\sigma^2}{\sigma^2 + \tau^2} \mu$$

$$- a_{\lim} A_{i\cdot} \frac{\tau^2}{\sigma^2 + \tau^2} \frac{1}{m} \sum_{h=1}^{m} y_h - a_{\lim} A_{i\cdot} \frac{\sigma^2}{\sigma^2 + \tau^2} \mu \le k \Bigg]$$

$$= P \Bigg[ A_{i\cdot} \frac{1}{m} \sum_{h=1}^{m} y_h + a_{\lim} \sum_{h=1}^{m} A_{ih} y_h - a_{\lim} A_{i\cdot} \frac{1}{m} \sum_{h=1}^{m} y_h$$

$$\le \frac{k - \frac{\sigma^2}{\sigma^2 + \tau^2} A_{i\cdot} \mu}{\frac{\tau^2}{\sigma^2 + \tau^2}} \Bigg] .$$

(2.43)

Furthermore,

$$A_{i\cdot} \frac{1}{m} \sum_{h=1}^{m} y_h + a_{\lim} \sum_{h=1}^{m} A_{ih} y_h - a_{\lim} A_{i\cdot} \frac{1}{m} \sum_{h=1}^{m} y_h$$

$$\sim N \Bigg[ A_{i\cdot} \mu, \left( \sigma^2 + \tau^2 \right) \frac{A_{i\cdot}^2 - a_{\lim}^2 A_{i\cdot}^2 + m a_{\lim}^2 \sum_{h=1}^{m} A_{ih}^2}{m} \Bigg] .$$

Using the expression for $a_{\lim}$ in (2.36),

$$E \left[ F_{\beta*} (k) - F_\beta^B (k) \right] = \frac{1}{n} \sum_{i=1}^{n} \left\{ P \Bigg[ z_i \le \frac{k - \sum_{h=1}^{m} A_{ih} \mu}{\tau \sqrt{\sum_{h=1}^{m} A_{ih}^2 - \frac{\sigma^2 \left( \sum_{h=1}^{m} A_{ih} \right)^2}{m(\sigma^2 + \tau^2)}}} \Bigg] \right.$$

$$\left. - P \Bigg[ z_i \le \frac{k - \sum_{h=1}^{m} A_{ih} \mu}{\tau \sqrt{\sum_{h=1}^{m} A_{ih}^2}} \Bigg] \right\} .$$

(2.44)

Figures 2.3 and 2.4 show the bias of $F_{\beta'}$ and $F_{\beta*}$, respectively, for different values of the ratio $\sigma^2 / \tau^2$ and assuming $\mu = 0$. The patterns of bias are similar to one another and to those seen in Figures 2.1 and 2.2: bias is high when the measurement error $\sigma^2$ is large relative to the prior variance $\tau^2$. In Figure 2.4, when the variance $(\tau^2)$ of the prior distribution of $\theta$ is large compared to the variance of the measurement error $(\sigma^2)$, $\sigma^2 \left( \sum_{h=1}^{m} A_{ih} \right)^2 / \left( m \left( \sigma^2 + \tau^2 \right) \right)$ is close to 0 and there is very little bias for all

values of $k$ as shown in the upper left panel. When the variance of the prior distribution of $\theta$ is small compared to the measurement error, $\sigma^2 \left(\sum_{h=1}^{m} A_{ih}\right)^2 / \left(m\left(\sigma^2 + \tau^2\right)\right)$ approaches $\left(\sum_{h=1}^{m} A_{ih}\right)^2 / m$ (for fixed $m$) and the effective scale factor of the first probability in (2.44) is less than $\tau \sqrt{\sum_{h=1}^{m} A_{ih}^2 - \frac{\sigma^2 \left(\sum_{h=1}^{m} A_{ih}\right)^2}{m(\sigma^2 + \tau^2)}}$ so the bias is large (in absolute value) for values of $k$ between the mean and the extrema as shown in the lower right panel. Note again that the bias of $F_{\beta*}$ is smaller than the bias of $F_{\beta'}$ by three orders of magnitude.

The arguments leading to expressions (2.41) and (2.44) provide heuristic evidence that $F_{\beta'}$ is biased for $F_{\beta}^B$ while $F_{\beta*}$ is asymptotically unbiased for $F_{\beta}^B$. We consider the case in which the elements of $\boldsymbol{A}$ are non-negative. This would be the case if $\boldsymbol{A}$ was a means or a sums matrix, i.e., a matrix in which the rows of $\boldsymbol{A}$ sum or average the $\theta_h$ within some larger grouping. A reasonable asymptotic assumption is that the row sums of $\boldsymbol{A}$ are bounded as $n \to \infty$. This holds, for example, if the number of small areas $n_i$ in group $i$ is bounded as the number of groups, $n$, goes to infinity. More specifically, let $\sum_{h=1}^{m} A_{ih} < c_1 < \infty$ and $\left(\sum_{h=1}^{m} A_{ih}\right)^2 < c_2 < \infty$ independent of $i$ as $m, n \to \infty$. Then we can say from equation (2.44) that the bias of $F_{\beta*}$ decreases as $m$ increases. Thus, preliminary analytical results in (2.41) and (2.44) plus numerical computations illustrated in Figures 2.3 and 2.4 suggest that $\boldsymbol{A}\theta^{CB}$ is in fact a better ensemble estimator than $\boldsymbol{A}\theta^{B}$.

Figure 2.3: Bias of $F_{\beta'}(k)$. The value of $\tau^2$ is fixed at 1 and $\sigma^2 =$ $.05, 0.1, 0.2, 0.5, 1, 2, 5, 10, 20$. $\boldsymbol{A}$ is an $1000 \times 200$ sums matrix where each $\beta_i$ is the sum of five $\theta_h$s. Here, we assume $\boldsymbol{\mu}$ is the 0 vector.

48



Figure 2.4: Bias of $F_{\theta CB}(k)$. The value of $\tau^2$ is fixed at 1 and $\sigma^2 = 0.05, 0.1, 0.2, 0.5, 1, 2, 5, 10, 20$. $A$ is an $1000 \times 200$ sums matrix where each $\beta_i$ is the sum of five $\theta_h$s. Here, we assume $\mu$ is the 0 vector.

Chapter 3

GENERAL CONSTRAINED BAYES WITH CAR

## 3.1   Introduction

As mentioned in Chapter 2, Stern and Cressie (1999) presented a method for

computing CB estimates for lattice data using a CAR model. In their method, the

spatial dependence parameter $\phi$ is fixed. In order to determine the best value of

$\phi$, they compute the value of the loss function for several different values of $\phi$ to

see which results in the smallest loss. They then demonstrate their method on the

Scotland lip cancer data set (Breslow and Clayton, 1993). In the first section of

this chapter, we try to reproduce their analysis using the GCB algorithm outlined

in Chapter 2.

We then apply the GCB algorithm to water quality data from the Mid-Atlantic

highlands, which is the original motivation for exploring constrained Bayes estimation

and the reason for the development of the GCB algorithm. The 2002 report titled

"Response of surface water chemistry to the Clean Air Act Amendments of 1990"

(Stoddard et al., 2002) was written to assess the response in surface water chemistry

in the northern and eastern United States to changes in acid deposition, primarily

from acid rain. A key indicator of whether or not a particular watershed is at risk

due to acid deposition is acid neutralizing capacity (ANC), the ability of a body of water to buffer inputs of acid. Stoddard et al. (2002) differentiate between acidic and acidified waters. 'Acidic' describes a condition that can be defined and documented, such as $ANC \leq 0$. 'Acidified' refers to an increase in acidity over time and does not require that the body of water be acidic. An important aspect of assessing the impact of the Clean Air Act Amendments of 1990 is identifying improvements in water chemistry. This speaks to quantifying and properly interpreting acidification, or more specifically, the change in ANC over time. Increasing values of ANC are expected in response to decreasing deposition of acid from the atmosphere. The region of interest consist of watersheds identified by the United States Geological Survey (USGS) and will be described in more detail later in this chapter. A spatially indexed estimator of ANC over time is needed for each watershed in the study area. Additionally, an ensemble estimate that can be used to identify correctly the proportion of watersheds with decreasing ANC will be valuable. The spatial structure of the watersheds from which the water quality data come is hierarchical or nested; each watershed on one level of the hierarchy is divided into smaller watersheds. This suggests that spatial correlation may exist at several levels of the hierarchy. We specify a CAR model on two levels of the hierarchy and demonstrate the ability of the GCB algorithm to handle such a model.

## 3.2 Application: lip cancer counts in Scotland

In the Scotland lip cancer data set, 56 geographic districts are defined in the local government structure prior to a 1995 reorganization (Stern and Cressie, 1999). The data set consists of observed cases of lip cancer (O); expected cases of lip cancer (E) based on variation in demographic factors across the region; the percent of the population in each district engaged in an outdoor industry, such as agriculture, fishing, or forestry (denoted AFF); and the neighbors (adjacent districts) for each district. The analysis by Stern and Cressie assumes the model

$$\boldsymbol{Y}|\boldsymbol{\theta} \sim N\left(\boldsymbol{\theta}, \sigma^2 \boldsymbol{D}\right)$$

$$\boldsymbol{\theta} \sim N\left(\boldsymbol{X}\boldsymbol{\gamma}, \tau^2 \left(\boldsymbol{I} - \phi\boldsymbol{C}\right)^{-1} \boldsymbol{V}\right), \tag{3.1}$$

where

$$\boldsymbol{\theta} = E\left(\boldsymbol{Y}\right)$$

and

$$Y_h = \left(\frac{O_h}{E_h}\right)^{1/2} + \left(\frac{O_h + 1}{E_h}\right)^{1/2}. \tag{3.2}$$

The vector $\boldsymbol{x}_h^T = (1\ x_h)$ consists of an intercept and $x_h =$ the AFF for district $h$. The other elements of the model are defined according to (iv)-(vi) in Stern and Cressie (1999) and are as follows:

- $\boldsymbol{D}$ is diagonal with $h$th element $d_{hh} = 1/E_h$.

- $\boldsymbol{V}$ is diagonal with $h$th element $d_{hh} = 1/E_h$.

- $c_{hk} = \begin{cases} (E_k/E_h)^{1/2}, & k \in N_h \\ 0, & \text{elsewhere}, h = 1, \ldots, 56, \end{cases}$

where $N_h$ is the neighborhood of district $h$. Additionally, $\sigma^2$ is fixed at 1. Stern and Cressie (1999) compute the *predicted mean squared error* (PMSE) for several values of $\phi$ and find that the value of $\phi$ that yields the smallest PMSE is 0.14.

We also use the model in (3.1) but consider $\phi$ random. The prior distributions on each of the parameters in our fully Bayesian algorithm are listed below in (3.3). The normal priors, $N(\mu, \tau)$, specify the mean $\mu$ and precision $\tau$ rather than the usual mean and variance. The priors are

$$
\begin{aligned}
\sigma^2 &\sim \text{Inverse gamma}\,(0.001, 0.001) \\
\gamma_i &\sim N\,(0, 0.0001)\,, \text{ for } i = 0, 1 \\
\tau^2 &\sim \text{Inverse gamma}\,(0.001, 0.001) \\
\phi &\sim \text{Uniform}\,(\phi_{\min}, \phi_{\max})\,,
\end{aligned}
\tag{3.3}
$$

where $\phi_{\min}$ and $\phi_{\max}$ are defined as in (1.11). These priors were chosen to reflect the fact that we have little information regarding the parameters ahead of time. The precision parameters $1/\sigma^2$ and $1/\tau^2$ have gamma distributions with mean 1 and variance 1000 which, on the scale of the data, represent uncertain precision. Likewise, the prior distributions on $\{\gamma\}$ suggest we have little information about the sign or magnitude of those parameters. The spatial correlation parameter $\phi$ is given equal density over its valid quantities.

Three separate MCMC chains were run with initial values given by:

- Chain 1

    - $\sigma^2 = 10000$

- $\gamma_0 = \gamma_1 = 0$

- $\tau^2 = 10000$

- $\phi = 0.1$

- Chain 2

  - $\sigma^2 = 10000$

  - $\gamma_0 = \gamma_1 = 10$

  - $\tau^2 = 10000$

  - $\phi = 0$

- Chain 3

  - $\sigma^2 = 100000$

  - $\gamma_0 = \gamma_1 = -10$

  - $\tau^2 = 100000$

  - $\phi = 0.5$

Each chain was run was run for 10000 iterations, the first 4000 of which were discarded as the "burn-in". Gelman-Rubin statistics (Gelman and Rubin, 1992) indicate convergence of all three chains. The posterior mean of $\phi$ was 0.1365 to four decimal places, agreeing quite well with the value of 0.14 obtained by Stern and Cressie (1999). The parameters of interest are the relative risks $\lambda_h = E(O_h)/E_h$. The relationship between $\theta$ and $\lambda$ is

$$\theta_h \approx 2\lambda_h^{1/2}\left(1 + \frac{1}{8\lambda_h E_h}\right), \tag{3.4}$$

the approximation due to a Taylor series expansion of $E(Y_h)$ around $E(O_h)$. Stern and Cressie (1999) compute the CB estimates of $\theta$ and then back-transform to obtain estimates of $\lambda$. Results from Section 2.4 demonstrate that the transformed CB estimates are not the same as the CB estimates of the transformed parameter. Thus, the loss function should be defined on $\lambda$ and its CB estimates be computed directly. However, to be consistent with Stern and Cressie (1999) we computed the GCB estimates for $\theta$ and then computed estimates for $\lambda$ using (3.4). The estimates of $\lambda$ from Stern and Cressie (1999) and $\lambda_{GCB}$ are shown in Table 3.1. With the exception of districts 17 and 55, the estimates are comparable. The transformation in (3.2) is such that $Y$ is approximately normal, and it generally preserves the ordering of the data: when the $\{O_h/E_h\}$ are placed in descending order, the $\{Y_h\}$ are generally in descending order, too. The values for $Y_{17}$ and $Y_{55}$ are conspicuously out of order. The two districts have six and five neighbors, respectively, which are about the average number of neighbors for a given district. Thus, it seems unlikely that that the inconsistent estimates for these two districts are due to lack of information from its "neighborhood".

## 3.3 Application: acid neutralizing capacity trends in the Mid-Atlantic Highlands

The problem of CB small area estimation with a general covariance matrix is now examined in an environmental setting. The response of interest is a trend parameter, namely the change in acid neutralizing capacity (ANC) over time. ANC measures

Table 3.1: Constrained Bayes (CB) estimates and General Constrained Bayes (GCB) estimates of the relative risk for skin cancer of 56 Scotland districts. The fixed $\phi$ estimate is the CB estimate from Stern and Cressie (1999) with $\phi = 0.14$; the random $\phi$ estimate is the GCB estimate using the method presented in Section 2.3.1. With the exception of districts 17 and 55, the estimates are similar.

| | CB:$\phi$ fixed | GCB:$\phi$ random | | CB:$\phi$ fixed | GCB:$\phi$ random |
|---|---|---|---|---|---|
| $\lambda_1$ | 6.77 | 6.71 | $\lambda_{29}$ | 1.25 | 1.15 |
| $\lambda_2$ | 4.22 | 4.49 | $\lambda_{30}$ | 1.04 | 1.12 |
| $\lambda_3$ | 3.42 | 3.78 | $\lambda_{31}$ | 1.07 | 1.14 |
| $\lambda_4$ | 3.64 | 3.78 | $\lambda_{32}$ | 1.36 | 1.22 |
| $\lambda_5$ | 3.58 | 3.57 | $\lambda_{33}$ | 1.05 | 1.07 |
| $\lambda_6$ | 3.58 | 3.53 | $\lambda_{34}$ | 0.88 | 0.99 |
| $\lambda_7$ | 3.15 | 3.24 | $\lambda_{35}$ | 0.90 | 0.93 |
| $\lambda_8$ | 2.90 | 3.21 | $\lambda_{36}$ | 0.85 | 0.93 |
| $\lambda_9$ | 3.32 | 3.22 | $\lambda_{37}$ | 0.93 | 0.91 |
| $\lambda_{10}$ | 3.08 | 3.09 | $\lambda_{38}$ | 0.74 | 0.89 |
| $\lambda_{11}$ | 3.01 | 3.04 | $\lambda_{39}$ | 1.01 | 0.91 |
| $\lambda_{12}$ | 3.29 | 3.04 | $\lambda_{40}$ | 0.71 | 0.84 |
| $\lambda_{13}$ | 3.53 | 3.16 | $\lambda_{41}$ | 0.52 | 0.56 |
| $\lambda_{14}$ | 2.50 | 2.57 | $\lambda_{42}$ | 0.65 | 0.55 |
| $\lambda_{15}$ | 2.04 | 2.22 | $\lambda_{43}$ | 0.66 | 0.59 |
| $\lambda_{16}$ | 2.24 | 2.07 | $\lambda_{44}$ | 0.43 | 0.44 |
| $\lambda_{17}$ | 2.74 | 2.26 | $\lambda_{45}$ | 0.41 | 0.39 |
| $\lambda_{18}$ | 1.58 | 1.76 | $\lambda_{46}$ | 0.45 | 0.43 |
| $\lambda_{19}$ | 1.84 | 1.72 | $\lambda_{47}$ | 0.39 | 0.44 |
| $\lambda_{20}$ | 1.58 | 1.69 | $\lambda_{48}$ | 0.37 | 0.38 |
| $\lambda_{21}$ | 1.48 | 1.56 | $\lambda_{49}$ | 0.34 | 0.32 |
| $\lambda_{22}$ | 1.46 | 1.39 | $\lambda_{50}$ | 0.39 | 0.34 |
| $\lambda_{23}$ | 1.30 | 1.31 | $\lambda_{51}$ | 0.38 | 0.43 |
| $\lambda_{24}$ | 1.14 | 1.32 | $\lambda_{52}$ | 0.36 | 0.40 |
| $\lambda_{25}$ | 1.20 | 1.25 | $\lambda_{53}$ | 0.24 | 0.26 |
| $\lambda_{26}$ | 1.10 | 1.23 | $\lambda_{54}$ | 0.22 | 0.21 |
| $\lambda_{27}$ | 1.13 | 1.24 | $\lambda_{55}$ | 0.16 | 0.07 |
| $\lambda_{28}$ | 1.14 | 1.16 | $\lambda_{56}$ | 0.16 | 0.15 |

how well bodies of water can buffer inputs of acidity from, for example, acid rain or industrial runoff. As such, it can be used to measure the risk of a body of water to acidification. The geographical region of interest is the Mid-Atlantic Highlands, an area roughly bounded by New York to the north, North Carolina to the south, and

Ohio to the west. This region is particularly susceptible to acid rain as a result of its location relative to industrial centers and prevailing weather patterns. The inferential goal again is two-fold. First, we would like to obtain good *individual* estimates of $\theta_h = \Delta\text{ANC}_h/\texttt{time}$, for the $h = 1, \ldots, m$ watersheds in this region, and second, we would like the spatially indexed ensemble $\{\theta_h^{est}\}_{h=1}^m$ to have similar characteristics to the spatially indexed true values $\{\theta_h\}_{h=1}^m$. This will enable us, for example, to estimate $P(\theta < 0)$, the proportion of watersheds for which ANC is decreasing over time. Following a Gaussian-Gaussian probability model, we model $\beta$ as a function of auxiliary information and place a Gaussian distribution on $\theta_{observed} \mid \theta_{true}$.

In this section we compute constrained Bayes estimates for the watershed-specific rates-of-change in ANC for the data set described above. The data originally consisted of 103 hydrologic units (HUCs) in the Mid-Atlantic highlands. The HUCs for which we are estimating $\beta$ are among the smallest watershed unit identified by the USGS (Seaber et al., 1987). Within each of the 103 HUCs for which data were available, ANC was recorded for at least one site in at least one year between 1993 and 1998. Since we are interested in $\theta = \Delta\text{ANC}_h/\texttt{time}$, 17 HUCs had to be discarded because ANC was available for only one year. Thus, we were left with 86 HUCs in which we had ANC data for at least two years. From these data, we calculated $\hat{\theta} = \theta^{OLS}$, the ordinary least squares slope when ANC was regressed on year. HUC level covariates that were available included spatial coordinates of the HUC centroid, area, average elevation, average topological slope, maximum slope, percents

agriculture, urban, and forest, and $SO_4$ deposition for the current and three previous years.

The HUC structure is nested. The United States is divided and sub-divided into successively smaller hydrologic units which are classified into four levels: regions, sub-regions, accounting units, and cataloging units, according to the USGS (http://water.usgs.gov/GIS/huc.html). The hydrologic units are arranged within each other, from the smallest (cataloging units) to the largest (regions). Each hydrologic unit is identified by a unique hydrologic unit code consisting of eight digits based on the four levels of classification in the hydrologic unit system. Thus, 8 digits are required to identify a unique cataloging unit, 6 digits are required to identify a unique accounting unit, 4 digits are required to identify a unique sub-region, and 2 digits are required to identify a unique region. When we refer to an 8-digit HUC, we are referring to a cataloging unit, and so forth for other levels in the hierarchy.

The HUC structure is not arbitrary. Major geographic regions are contained within 2-digit HUCs; large river systems and coastal drainage basins are contained within 4-digit HUCs; 6-digit HUCs contain surface drainage basins or combination of basins; and 8-digit HUCs distinguish parts of drainage basins and unique hydrologic features. It is thus not unreasonable to assume that all of the 8-digit HUCs nested within the same 6-digit HUC have similar hydrological characteristics and likewise for other nesting structures in the hierarchy. It is under this assumption that we justify the use of the CAR model and create the spatial adjacency matrix $C$.

### 3.3.1 GCB estimates for a one-stage CAR

We define the matrix $C$ so that it reflects appropriate relationships among the different 8-digit HUCs; this model assumes that all (or at least a majority) of the spatial dependence resides among the 8-digit HUCs nested within the same 6-digit HUCs. Hence, we assume the model

$$
\begin{aligned}
\hat{\theta}_h &= \theta_h + \epsilon_h \\
\theta_h &= x_h^T \gamma + \omega_h,
\end{aligned}
\tag{3.5}
$$

where

$$
\epsilon \sim N\left(0, \sigma^2 D\right)
\tag{3.6}
$$

and

$$
\omega \sim N\left(0, \tau^2 \left(I - \phi C\right)^{-1} M\right).
\tag{3.7}
$$

Additionally, we will define $\left(I - \phi C\right)^{-1} M \equiv \Phi$.

First, we define the matrices $D$, $C$, and $M$ from (3.6) and (3.7). The matrix $D$ defines the covariance structure among the $\hat{\theta}$ vector. If we think of $\hat{\theta}_h$ as the true value $\theta_h$ plus some observation or measurement error, it is reasonable to assume that the precision of $\hat{\theta}_h$ is independent of $\hat{\theta}_k$ for $k \neq h$, and is proportional to $n_h$, where $n_h$ is the number of observations taken in HUC $h$ and used in the computation of $\hat{\theta}_h$. Therefore, we set $D = \text{diag}\left(n_1^{-1}, n_2^{-1}, \ldots, n_m^{-1}\right)$.

The matrix $C$ defines the existence and strength of association among the different HUCs. The neighborhood structure we use reflects the nested structure of the USGS HUC designation. All HUC-8s within the same HUC-6 region are considered

part of the same neighborhood. No spatial relationships among HUC-6 or HUC-4 regions are considered at this time. The adjacency matrix $C$ is then defined as

$$
c_{hk} = \begin{cases} \dfrac{1}{\text{distance between HUC-8 centroids}}, & \text{if } h \text{ and } k \text{ are neighbors}; \\[2em] 0, & \text{otherwise}. \end{cases} \tag{3.8}
$$

An adjacency matrix defined this way indicates that HUC-8 regions within the same HUC-4 are correlated, and the strength of the correlation is inversely proportional to the distance between the HUC-8 centroids.

The matrix $M$ must be a diagonal matrix (Cressie, 1993) and each element of the diagonal is a scaling factor for the variance in each 8-digit HUC. We assume that the more neighbors upon which a given HUC is dependent, the more of a stabilizing effect this will have on the response for that HUC. Another way to look at it is that, referring back to equations (1.6) and (1.7), the form of the conditional mean of $\theta_h | \boldsymbol{\theta}_{-h}$ includes a weighted average of the vector $\boldsymbol{\theta}_{-h}$. It is therefore not unreasonable to expect an averaging effect causing the conditional variance to decrease as the number of neighbors, and hence the number of terms in the conditional mean, increases. Thus we define $m_{hh} = N_h^{-1}$, where $N_h$ is the number of neighbors for HUC $h$. This choice for $M$ along with the choice for $C$ in (3.8) satisfies $c_{hk} m_{kk} = c_{kh} m_{hh}$ since for $h, k$ in the same neighborhood, $c_{hk} = c_{kh}$ and $m_{hh} = m_{kk}$; for $h, k$ in different neighborhoods, $c_{hk} = c_{kh} = 0$.

We place the following priors on the parameters in the model:

$$
\sigma^2 \sim \text{Inverse gamma}\,(0.001, 0.001)
$$

$$\gamma_h \quad \sim \quad N\,(0,1000)$$

$$\tau^2 \quad \sim \quad \text{Inverse gamma}\,(0.001,0.001)$$

$$\phi \quad \sim \quad \text{Uniform}\,(0,\phi_{\max})$$

where $\phi_{\max}$ is defined in (1.11). These priors are chosen because they are fairly uninformative in the sense that the variances are large. The precisions, $1/\sigma^2$, of $\hat{\theta}_h$, $h = 1,\ldots,m$ and $1/\tau^2$, of $\omega_h$, $h = 1,\ldots,m$ are gamma distributions with variances$= 0.001/0.001^2 = 1000$. Note that $\phi \in (0,\phi_{\max})$ leads to a valid covariance matrix since $\phi_{\min} < 0$. Results from fitting this model are described in Section 3.3.3.

### 3.3.2   GCB estimates for a two-stage CAR

In the previous section, we assumed all of the spatial association was among 8-digit HUCs within the same 6-digit HUC. By the nature of the HUC designation system, it is certainly plausible that all 6-digit HUCs within the same 4-digit HUC have similar ANC characteristics, and so on. We will redefine the model (3.5) to include two levels of spatial correlation, one among 8-digit HUCs within the same 6-digit HUC and another among 6-digit HUCs within the same 4-digit HUC. We cannot explore the existence of correlation among the 4-digit HUCs within the same 2-digit HUCs as all of the 4-digit HUCs in the region of study belong to the same 2-digit HUC.

To account for this second possible source of spatial correlation, we can expand the model to, say

$$\hat{\theta}_{ih} \quad = \quad \theta_{ih} + \epsilon_{ih}$$

$$\theta_{ih} = x_h^T \gamma + \alpha_i + \omega_h. \tag{3.9}$$

Here, $i$ indexes the 6-digit HUC and $h$ indexes the 8-digit HUC. Thus, $\alpha_i$ is a random effect for the 8-digit HUC and $\omega_j$ is a random effect for the 8-digit HUC. A CAR model can then be placed on both $\alpha$ and $\omega$ to account for spatial dependence that may occur on the larger scale. When the variation was additive as in (3.9), the algorithm had a difficult time partitioning it between the 8-digit HUCs and 6-digit HUCs. Rather than define the random effects as independent and additive, we nested them:

$$\hat{\theta}_{ih} = \theta_{ih} + \epsilon_{ih}$$
$$\theta_{ih} = x_h^T \gamma + \alpha_i + \omega_{h|i}, \tag{3.10}$$

where

$$\epsilon \sim N\left(0, \sigma^2 D\right), \tag{3.11}$$

$$\alpha \sim N\left(0, \tau_6^2 \left(I - \phi_6 C_6\right)^{-1} M_4\right), \text{ and} \tag{3.12}$$

$$\omega | \alpha \sim N\left(\alpha, \tau_8^2 \left(I - \phi_8 C_8\right)^{-1} M_8\right). \tag{3.13}$$

Again, $i$ indexes the 6-digit HUC and $h$ indexes the 8-digit HUC.

The matrices $D$, $C_8$, and $M_8$ are defined as in Section 3.3.1. $C_6$ and $M_6$ are defined similarly, so that

$$c_{6ij} = \begin{cases} \dfrac{1}{\text{distance between HUC-6 centroids}}, & \text{if } i \text{ and } j \text{ are neighbors;} \\ \\ 0, & \text{otherwise,} \end{cases} \tag{3.14}$$

and $M_6 = [m_{6ii}] = [N_{6i}^{-1}]$, $i = 1, \ldots, n_6$, where $N_{6i}$ is the number of 6-digit HUC neighbors for 6-digit HUC $i$.

The following priors are placed on the parameters in the model:

$$\sigma^2 \sim \text{Inverse gamma}(0.001, 0.001)$$

$$\gamma_h \sim N(0, 1000)$$

$$\tau_8^2 \sim \text{Inverse gamma}(0.001, 0.001)$$

$$\tau_6^2 \sim \text{Inverse gamma}(0.001, 0.001)$$

$$\phi_8 \sim \text{Uniform}(0, \phi_{8\max})$$

$$\phi_6 \sim \text{Uniform}(0, \phi_{6\max})$$

where $\phi_{8\max} = \eta_{8\max}^{-1}$, $\phi_{6\max} = \eta_{6\max}^{-1}$, and $\eta_{8\max}$ and $\eta_{6\max}$ are the largest eigenvalues of $M_8^{1/2} C_8 M_8^{1/2}$ and $M_6^{1/2} C_6 M_6^{1/2}$, respectively.

### 3.3.3 Results

Three independent MCMC chains were run for both the one-stage and two-stage CAR models, one with high starting values for the parameters, one with intermediate values, and one with low values. Convergence was reached on all chains for all parameters after a burn-in of 10000 iterations according to Gelman-Rubin statistics (see Table 3.2 for a summary). The chains were then run for an additional 10000 iterations which were used in the analysis.

The posterior means of some of the parameters in the model are shown in Table 3.3. Of particular note are the estimates of $\sigma^2$ and $\tau^2$. Recall from the model

Table 3.2: Gelman-Rubin ratios for the parameters in the one-stage model. Rather than list the statistic for each $\beta$ and $\gamma$, we list only the maximum value.

| Parameter | G-R Statistic |
|---|---|
| $\theta$ | 1.002 (max) |
| $\gamma$ | 1.001 (max) |
| $\phi$ | 1.000 |
| $\sigma^2$ | 0.998 |
| $\tau^2$ | 0.969 |

(2.1) that the parameter $\sigma^2$ can be interpreted as measurement or observation error while $\tau^2$ represents a scaling factor on the variability in the spatial dependence. The results shown in Table 3.3 suggest that the variability in the data is primarily from the former source. The consequences of this become evident as we analyze the data further. The complete conditional distribution of $\boldsymbol{\theta}$, given the data and all of the other parameters is given by (2.18). In the Gibbs sampler, if $\sigma^2$ is large and $\tau^2$ is small, the approximate conditional distribution of $\boldsymbol{\theta}$, given the data and all of the other parameters is

$$N\left(\boldsymbol{X}\gamma, \boldsymbol{\Sigma}\right)$$

in the notation of (2.18), or using the notation of (3.5) through (3.7),

$$N\left(\boldsymbol{X}\gamma, \boldsymbol{\Phi}\right). \tag{3.15}$$

Thus, the realizations of $\boldsymbol{\theta}$ in the Gibbs sampler are essentially drawn from (3.15), and hence

$$E\left[\boldsymbol{\theta}^T\left(\boldsymbol{I}-\boldsymbol{P}\right)^T\boldsymbol{\Phi}^{-1}\left(\boldsymbol{I}-\boldsymbol{P}\right)\boldsymbol{\theta}\,|\,\hat{\boldsymbol{\theta}}\right]$$
$$\approx\ E\left(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}\right)^T\left(\boldsymbol{I}-\boldsymbol{P}\right)^T\boldsymbol{\Phi}^{-1}\left(\boldsymbol{I}-\boldsymbol{P}\right)E\left(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}\right). \tag{3.16}$$

Equation (3.16) says that, given the data and all other parameters, the posterior

Table 3.3: Posterior means of 10000 realizations for the spatial dependence parameter and the two variance components.

| One-stage CAR | | Two-stage CAR | |
|---|---|---|---|
| Parameter | Posterior mean | Parameter | Posterior mean |
| $\phi$ | 4.67 | $\phi_8$ | 0.24 |
| | | $\phi_6$ | 3.91 |
| $\tau^2$ | 1.67 | $\tau_8^2$ | 21.96 |
| | | $\tau_6^2$ | 1.25 |
| $\sigma^2$ | 793165.10 | $\sigma^2$ | 161810.80 |

mean of the residuals from the projection of $\boldsymbol{\theta}$ onto the column space of $\boldsymbol{X}$ is approximately equal to the residuals from the projection of the Bayes estimators onto the column space of $\boldsymbol{X}$. Equation (3.16) is also an approximation of (2.3) with $t\left(\boldsymbol{Y}\right) = E\left(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}\right)$, implying that the GCB estimates are approximately equal to the Bayes estimates. This result is illustrated in Figure 3.1 which shows the empirical distribution functions (edf) of the Bayes estimates of $\boldsymbol{\theta}$ (dashed/dotted line) and of the GCB estimates of $\boldsymbol{\theta}$ (dotted line). The solid line represents the posterior mean of $F_\theta$, the edf of $\boldsymbol{\theta}$, which is the distribution we are trying to match. The edf of the CB estimates follows the posterior mean of $F_\theta$, as we would expect. That the edf of the Bayes estimates also follows the posterior mean of $F_\theta$ and is difficult to distinguish from the edf of the CB estimates is because the estimate of $\tau^2$ is so small compared to $\sigma^2$. Thus, there is very little shrinkage in the Bayes estimates.

To assess the amount of spatial correlation within the different levels of HUC hierarchy, Moran's I was calculated on two different scales, first to measure the dependence among 8-digit HUCs within the same 6-digit HUC, and second to measure the dependence among the 6-digit HUCs within the same 4-digit HUCs. There were not enough data to compute Moran's I at other scales. Results from Moran's I indi-

Figure 3.1: Empirical distribution functions.

cate that there is no spatial dependence at the 8-digit HUC level ($p = 0.94$), which is consistent with the results above, while at the 6-digit HUC level, there appears to be very minimal spatial dependence ($p = 0.08$). Figures 3.2(a) and 3.2(b) show contour plots of the ANC rates predicted by GCB with the one-stage CAR and two-stage CAR, respectively. The patterns are similar. However, the image in Figure 3.2(b) is bit smoother, which reflects the slight correlation among the 6-digit HUCs. While the results of this comparison are largely negative, the main point is to illustrate the possibility of computing GCB estimates in fairly complex models, motivated by a real application.

(a) One-Stage CAR



(b) Two-Stage CAR

Figure 3.2: Contour plots of GCB estimated ANC rates in $\mu$g/L/year. Darker shades represent lower levels. The black contour line is the 0 contour

Chapter 4

# SMALL AREA ESTIMATION WITH GENERAL CONSTRAINED

# BAYES USING PENALIZED SPLINES

## 4.1 Introduction

In this chapter, the GCB algorithm from Chapter 2 for computing CB estimates with a general covariance matrix is applied to a spatial small area model specified with penalized radial basis functions and formulated as a linear mixed model as described in Chapter 1. We analyze the data set from Opsomer et al. (2008). The data are from the Environmental Monitoring and Assessment Program (EMAP) of the U.S. Environmental Protection Agency in which 334 lakes out of 21,026 in the Northeastern U.S. were surveyed over six years from 1991 to 1996. There are 551 observations in the data set. We will again estimate ANC for each 8-digit HUC; however, whereas in Chapter 3 we modeled change in ANC over time for each HUC, here we will follow the analysis of Opsomer et al. (2008) and model the mean ANC for each HUC. The model used by Opsomer et al. (2008) is of the form

$$y_{hi} = f(x_{hi}) + \epsilon_{hi} \qquad (4.1)$$

where $y_{hi}$ is ANC for the $i$th observation in the $h$th HUC. Also, $f$ here is defined as

$$f\left(x_{hi}\right) = \beta_0 + \beta_1 x_{hi} + u_h + \sum_{k=1}^{K} \beta_{k+1} r_{k+1}\left(c_{ti}\right), \qquad (4.2)$$

where $x_{hi}$ is the elevation at the $i$th observation in HUC $h$, $u_h \sim N\left(0, \sigma_u^2\right)$ is a random effect for HUC $h$, $r_k\left(c_{hi}\right)$ is the *transformed radial basis* function defined in Ruppert et al. (2003, p. 253) as

$$r_{k+1}\left(c_{hi}\right) = \sum_{k'=1}^{K} C\left(c_{hi} - \kappa_k\right)\left[C\left(\kappa_k - \kappa_{k'}\right)\right]^{-1/2}, \qquad (4.3)$$

$C\left(r\right) = \|r\|^2 \log\|r\|$, $c_{hi} = \left(c_{1hi}, c_{2hi}\right)$ denotes the geographical coordinates for observation $i$ in HUC $h$, and $\kappa_k$, $k = 1, \ldots, K$ are the geographical coordinates at the spline knots. Following Ruppert et al. (2003) as outlined in Chapter 1, (4.2) can be written as

$$\boldsymbol{Y} = \boldsymbol{X\beta} + \boldsymbol{Du} + \boldsymbol{Z\gamma} + \boldsymbol{\epsilon}, \qquad (4.4)$$

where $\boldsymbol{X}$, $\boldsymbol{D}$, and $\boldsymbol{Z}$ are the design matrices for the fixed effects, random HUC effect, and random knot effects, respectively, and

$$\boldsymbol{u} \sim \left(\boldsymbol{0}, \sigma_u^2 \boldsymbol{I}\right)$$

$$\boldsymbol{\gamma} \sim \left(\boldsymbol{0}, \sigma_\gamma^2 \boldsymbol{I}\right)$$

$$\boldsymbol{\epsilon} \sim \left(\boldsymbol{0}, \sigma_\epsilon^2 \boldsymbol{I}\right). \qquad (4.5)$$

The model in (4.4) is defined on each observation, but the characteristic Opsomer et al. (2008) wish to estimate is mean ANC

$$\bar{y}_h = \bar{\boldsymbol{x}}_h \boldsymbol{\beta} + u_h + \bar{\boldsymbol{z}}_h \boldsymbol{\gamma} \qquad (4.6)$$

for HUC $h$, where $\bar{x}_t$ is the mean elevation of HUC $h$ and $\bar{z}_h$ is the mean value of (4.3) for HUC $h$. As an estimator for $\bar{y}_h$, Opsomer et al. (2008) use

$$\hat{y}_h = \bar{x}_h \hat{\beta} + \hat{u}_h + \bar{z}_h \hat{\gamma}, \qquad (4.7)$$

where

$$
\begin{aligned}
\hat{\beta} &= \left(X'V^{-1}X\right)^{-1} X'V^{-1}Y \\
\hat{u} &= \sigma_u^2 D'V^{-1}\left(Y - X\hat{\beta}\right) \\
\hat{\gamma} &= \sigma_\gamma^2 Z'V^{-1}\left(Y - X\hat{\beta}\right) \\
V &= \sigma_u^2 DD' + \sigma_\gamma^2 ZZ' + \sigma_\epsilon^2 I, \qquad (4.8)
\end{aligned}
$$

and variance components are set equal to their *restricted maximum likelihood* (REML) estimates.

## 4.2 Estimation of mean ANC using GCB

We estimate $\bar{y}_h$ directly through (4.6) with the following model specification:

$$
\begin{aligned}
\bar{y}|\beta,\gamma,u &\sim N\left(\bar{x}\beta + u + \bar{z}\gamma, \sigma_\epsilon^2 \text{diag}\left(n_h\right)\right) \\
\beta &\sim N\left(0, \sigma_\beta^2 I\right) \\
u &\sim N\left(0, \sigma_u^2 I\right) \\
\gamma &\sim N\left(0, \sigma_\gamma^2 I\right), \qquad (4.9)
\end{aligned}
$$

where $n_h$ is the number of observations in HUC $h$. Results from Opsomer et al. (2008) suggest that the variances are on the order of $10^3$ to $10^5$. Therefore, all variances in (4.9) have a prior distribution of inverse gamma (0.5, 0.005). With this prior on $\sigma_\epsilon^2$,

Table 4.1: Analysis of the Northern Lakes data. Comparison of posterior means for the HUC-level model (based on 3000 Gibbs draws after burn-in) and ML/REML estimates for the site-level model in Opsomer et al. (2008).

| Parameter | Posterior mean | Opsomer et al. |
|---|---|---|
| Intercept | 169.1 | 228.6 |
| Elevation | -0.6633 | -0.814 |
| $\sigma_u$ | 49.8 | 71.2 |
| $\sigma_\gamma$ | 185.7 | 365.7 |
| $\sigma_\epsilon$ | 166.4 | 179.5 |

$\sigma_\beta^2$, $\sigma_u^2$, and $\sigma_\gamma^2$, the prior precisions for the distributions of $\bar{y}|\beta,\gamma,u$; $\beta$; $u$; and $\gamma$ have mean 100 and variance 20000. Such a prior is relatively non-informative.

Note that model (4.9) is different from the model in (4.4)–(4.5). This is necessary in order to try to duplicate the results from Opsomer et al. (2008). We do not want to estimate site responses and average them for each HUC as is done in Opsomer et al. (2008) since this would mean computing a transformation of GCB estimates. As already shown, this is not the same as computing the GCB estimate of the transformation. We expect the GCB estimates of $\{\bar{y}\}$ from (4.9) to be comparable to the estimates of $\{\bar{y}\}$ from (4.4)–(4.5). We do not necessarily expect the parameter estimates to be comparable.

## 4.3 Results

We again used WinBUGS (Lunn et al., 2000) to generate realizations from the conditional distributions. Three independent MCMC chains were run, one with high starting values for the parameters, one with intermediate values, and one with low values. Convergence is reached on all chains, based on Gelman-Rubin statistics, after a burn-in of 2000 iterations. The chains are then run for an additional 3000

iterations which are used in the analysis. Bayes estimates for the fixed and random effects are shown in Table 4.1. The Bayes estimates of the parameters are of the same order of magnitude as those from Opsomer et al. (2008), but they compare only roughly in value. We can argue that since we are modeling means rather than individual observations, we might expect the variance component for the spline $\sigma_\gamma^2$ to be smaller for model (4.9) than for model (4.4)–(4.5). However, we also observe that the estimate of the variance component for the HUC random effect $\sigma_u^2$ is smaller for model (4.9) than for model (4.4). This makes less sense since there is no averaging of the HUC effect. The estimates of the error variance $\sigma_\epsilon^2$ are much closer to one another compared to the estimates of the other variance components. This is likely due to the fact that in (4.9), the conditional variance of $\bar{y}$ is explicitly scaled in the model. The estimate of the elevation effect has the same sign and comparable magnitude in both models, which is what we would want.

The empirical distribution functions of the observed mean ANC, the Bayes estimates of mean ANC, and the GCB estimates of mean ANC for each HUC are shown in Figure 4.1. The edf of the Bayes estimates is shown by the thick gray line and clearly shows significant shrinkage. The edf of the GCB estimates, shown by the thick black line, corrects for this and follows much more closely the edf of the observed data, shown by the thin black line. For the most part, the GCB estimates preserve the ordering observed in the data as shown in Figure 4.2. For example, the minimum observed mean ANC is $6.95\mu g/L$ observed in HUC 2040301. The GCB estimate for mean ANC in HUC 2040301 is $-1.70\mu g/L$ and is the smallest of the

Figure 4.1: Empirical distribution functions of observed and estimated ANC Hydrologic Unit Code (HUC) small areas within Northeast U.S. The edfs indicate much less variability among the Bayes estimates than among either the observed data or the GCB estimates.

GCB estimates. The observed HUC means are plotted along the upper horizontal line and the GCB estimates are plotted along the lower horizontal line. Segments connect the observed value and GCB estimate for each HUC.

Finally, we compare the results from the GCB algorithm to the results from Opsomer et al. (2008). In Figure 4.3, each rectangle is located roughly at the centroid of a HUC for which ANC was observed in at least one site. More specifically, the latitude and longitude coordinates at sampled sites were averaged within a HUC; the rectangles are centered at these coordinates. Break points in the figure's gray scale for estimated mean ANC are the same as in Figure 4 of Opsomer et al. (2008).

Figure 4.2: Comparison of observed mean ANC ($\mu$g/L) and GCB estimates. Segments connect the observed mean ANC and GCB estimated mean ANC for each HUC. In general the GCB estimates do a good job of preserving the observed order of ANC values. Several of the segments in the middle of the plot cross, indicating that the ordering is not completely preserved.

Figure 4.3 indicates that mean ANC is low in northern coastal regions and in the higher elevations of the Adirondack Mountains (northern New York) and in the the White Mountains (New Hampshire and western Maine). High mean ANC is seen in the Great Lakes area of New York and in the Connecticut River valley. These are the same findings as in Opsomer et al. (2008).

Figure 4.3: GCB estimated ANC ($\mu$g/L) for Hydrologic Unit Code (HUC) small areas within Northeast U.S. Rectangles are centered at the mean latitude and longitude of observed sites within HUCs.

# Chapter 5

# TWO-STAGE, MODEL-ASSISTED ESTIMATION USING PENALIZED SPLINES

## 5.1 Introduction

In a survey framework, a two-stage sampling design can be employed to make the best use of what are often limited time and financial resources. Even with the ability to focus such resources, it is often the case that the sample sizes are not sufficiently large to make model-free inferences. The presence of auxiliary information suggests employing a model in our inferences. Opsomer et al. (2008) propose incorporating this auxiliary information through a class of model-assisted estimators based on penalized spline regression in single stage sampling. Zheng and Little (2003) also use penalized spline regression in a model-based approach for finite population estimation in a two-stage sample. In a survey context, weights computed from a set of auxiliary information are often applied to many study variables. With this approach, model-assisted estimators should fare better than model-based estimators as discussed in Section 5.2.3. We compare the two through a series of simulations.

76

If the weights computed from the auxiliary information are obtained with simple linear regression, applying them to many study variables does not pose much of a problem as the hat matrix is a function of the predictors only. However, when we replace a simple linear model with something more flexible such as a penalized spline, obtaining an adequate set of survey weights becomes more complicated. The hat matrix is now a function of the variance components associated with the knots and with any additional random effects assumed in the model. To address this, we fix the degrees of freedom of the smooth and estimate the variance components of the random effects assuming a linear model. Asymptotic properties of these estimators are examined both analytically and through simulation. It is shown that model-assisted estimators fare at least as well as model-based estimators in this context.

When we consider two-stage sampling with auxiliary information, we must specify the extent to which the auxiliary information is known and on what scale it is known. Four different cases are possible:

**Case A**: *The auxiliary information is available for all clusters in the population*

- Leads to regression modeling of quantities associated with the clusters, such as cluster totals

- Cluster quantities can be computed for all clusters

- Population quantities can be computed from cluster estimates

**Case B**: *Complete Element Level Auxiliaries*

- The auxiliary information is available for all elements in the population

- Leads to regression modeling of quantities associated with the elements

- Cluster and population quantities can then be computed from element estimates and observations

**Case C**: *Limited Element Level Auxiliaries*

- The auxiliary information is available for all elements in selected clusters only

- Leads to regression modeling of quantities associated with the elements

- Regression estimators can be used for cluster-level quantities only for the clusters selected in the first-stage sample

**Case D**: *Limited Cluster Level Auxiliaries*

- The auxiliary information is available for all clusters in the first-stage sample

- Design-based estimator can be used for population quantities

- In some cases, good estimators for population quantities are not available

In this chapter, we focus on Case A in which the auxiliary information is known for all clusters in the population.

## 5.2 Background

### 5.2.1 The Finite Population and Superpopulation Concept

We begin by assuming a population of elements $U = \{1, \ldots, k, \ldots, N\}$. These elements are partitioned into clusters or primary sampling units (PSUs),

$U_1, \ldots, U_i, \ldots, U_{N_I}$, which are mutually exclusive and collectively exhaustive. Thus, we have

$$U = \bigcup_{i \in U_I} U_i \text{ and } N = \sum_{i \in U_I} N_i,$$

where $N_i$ is the number of elements or secondary sampling units (SSUs) in $U_i$ and $U_I = \{1, 2, \ldots, N_I\}$. Associated with each cluster $i \in U_I$ is a known vector of auxiliary information $\boldsymbol{x}_i = (x_{i1}, x_{i2}, \ldots, x_{ip})^T$. Associated with each element $k \in U$ is an unknown quantity of interest $y_k$. We focus on finding an estimator of

$$t_y = \sum_{k \in U} y_k = \sum_{i \in U_I} \sum_{k \in U_i} y_k = \sum_{i \in U_I} t_{yi},$$

where $t_y$ is the total for the quantity $y$ over the entire population, and $t_{yi}$ is the total for the quantity $y$ for PSU $i$.

When dealing with finite populations, it is useful to consider a *superpopulation model*, say $\xi$, which specifies a relationship between the auxiliary variables and the response of interest. Cassel et al. (1977) give several possible interpretations for the superpopulation concept:

1. The finite population is actually drawn from a larger universe.

2. The distribution $\xi$ is modeled to describe a random mechanism or process in the real world.

3. The distribution $\xi$ is considered as a prior distribution reflecting subjective belief, as in a Bayesian approach. The unobserved numbers among $y_1, \ldots, y_N$ may be looked upon as parameters for which we seek the posterior distribution,

given the sample.

4. The distribution $\xi$, while being associated neither with a process in the real world nor with an expression of subjective belief, is used simply as a mathematical device to make explicit the theoretical derivations. For example, one may be interested in knowing the various model formulations which justify the use of an intuitively appealing estimator such as the sample mean.

5. The superpopulation approach may be a useful device for incorporating the treatment of nonsampling errors in survey sampling.

## 5.2.2   Two-stage Sampling

Regardless of how we interpret the superpopulation, the assumption is that the finite population is a realization from the superpopulation, and our sample is obtained from this realization through some specified sampling design. We will use a probability sample. Probability sampling is an approach to sample selection that satisfies certain conditions, as outlined in Särndal et al. (1992):

1. We can define the set of samples, $\mathcal{S} = \{s_1, s_2, \ldots, s_M\}$ that are possible to obtain with the sampling procedure.

2. A known probability of selection $p(s)$ is associated with each possible sample $s$.

3. The procedure gives every element in the population a nonzero probability of selection.

4. We select one sample by a random mechanism under which each possible $s$ receives exactly the probability $p(s)$.

In the two-stage sampling design, we first select a sample of clusters, $s_I$ (the subscript $I$ is a Roman numeral referring to stage I), from the universe of clusters, $U_I$, based on some design $p_I(\cdot)$. The design $p_I(\cdot)$ assigns a *first order inclusion probability* $\pi_{Ii} > 0$, for $i \in U_I$. The $\pi_{Ii}$ is the probability that cluster $i$ is selected in the first stage of the sample. Additionally, $p_I(\cdot)$ assigns a *second order inclusion probability* $\pi_{Iij}$ to each pair of clusters $i, j$ in $U_I$. The $\pi_{Iij}$ is the probability that both clusters $i$ and $j$ are selected in the first stage. Note that $\pi_{Iii} = \pi_{Ii}$. In the second stage, from every cluster $i \in s_I$, a sample $s_i$ is drawn from $U_i$ according to a design $p_{II}(\cdot|s_I)$ with first and second order inclusion probabilities $\pi_{k|i}$ and $\pi_{kl|i}$, respectively.

### 5.2.3 Design Based Inference versus Model Based Inference

In the course of this chapter, we will look at *model-based* and *model-assisted* estimators. The model-based and model-assisted approaches view the role of the

model differently. In model-based estimation, the statistical properties of the estimator (expected value, bias, variance, etc.) are derived from the assumptions in the model; the sampling design is completely ignored although Särndal et al. (1992) advocate randomized selection of the sample as a safeguard against selection bias. Thus, we view the model not only as a random process in the real world, but we also use it to describe nonsampling error. In model-assisted estimation, the statistical properties of the estimator (expected value, bias, variance, etc.) are derived solely from the design $p(s)$. The model $\xi$ is simply a mechanism to describe how the finite population we want to study arises. Additionally, by explicitly stating the model, we can derive and justify the estimates for model parameters. However, the error properties of the model properties are ignored.

## 5.3    Model-assisted estimator with penalized spline

We focus on Case A, in which the auxiliary information is the scalar $x_i$ at the cluster level and is known for all clusters in the population. Consider the superpopulation model

$$\xi: \ y_k = f(x_i) + \alpha_i + \epsilon_k; \quad \text{for } k \in U_i \text{ and } i = 1, 2, \ldots, N_I, \tag{5.1}$$

where $\{\alpha_i\}$ are independent random variables with mean zero and variance $\tau^2$, $\{\epsilon_k\}$ are independent with mean zero and variance $\sigma^2$, and the two sequences are independent of one another.

Recall that the finite population is considered a realization from this model. If the entire realization had been observed, we could estimate $f(\cdot)$ using a P-spline

with truncated power basis, according to the details in Chapter 1, as

$$f\left(x;\boldsymbol{\beta}\right) = \beta_0 + \beta_1 x + \ldots + \beta_p x^p + \sum_{l=1}^{K} \beta_{p+l} \left[\left(x - \kappa_l\right)_+\right]^p. \tag{5.2}$$

Once again, $p$ is the degree of the spline, $\kappa_1 < \cdots < \kappa_K$ is a set of fixed knots, $\boldsymbol{\beta} = \left(\beta_0, \ldots, \beta_{p+K}\right)^T$ is the coefficient vector, and $(x)_+ = x I_{\{x>0\}}$. Substituting (5.2) into (5.1), we get

$$y_k = \beta_0 + \beta_1 x_i + \sum_{l=1}^{K} \beta_{1+l} \left(x_i - \kappa_l\right)_+ + \alpha_i + \epsilon_{ik}; \quad \text{for } k \in U_i \text{ and } i = 1, 2, \ldots, N_I. \tag{5.3}$$

For simplicity, we consider the linear case with independent errors, although the methods here can be extended to higher order polynomials and models with heteroscedasticity. In matrix notation, we can write this as

$$\boldsymbol{Y} = \boldsymbol{Z}_1 \boldsymbol{\beta}_1 + \boldsymbol{Z}_2 \boldsymbol{\beta}_2 + \boldsymbol{Z}_3 \boldsymbol{\beta}_3 + \boldsymbol{\epsilon}, \tag{5.4}$$

where

$$\boldsymbol{Y}_{N \times 1} = \begin{bmatrix} [y_k]_{k \in U_1} \\ [y_k]_{k \in U_2} \\ \vdots \\ [y_k]_{k \in U_{N_I}} \end{bmatrix}, \quad \boldsymbol{Z}_1 = \begin{bmatrix} [1 \ x_1]_{k \in U_1} \\ [1 \ x_2]_{k \in U_2} \\ \vdots \\ [1 \ x_{N_I}]_{k \in U_{N_I}} \end{bmatrix}_{N \times 2}, \quad \boldsymbol{\beta}_1 = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix},$$

$$\boldsymbol{Z}_2 = \begin{bmatrix} \left[(x_1 - \kappa_1)_+ \ (x_1 - \kappa_2)_+ \ \cdots \ (x_1 - \kappa_K)_+\right]_{k \in U_1} \\ \left[(x_2 - \kappa_1)_+ \ (x_2 - \kappa_2)_+ \ \cdots \ (x_2 - \kappa_K)_+\right]_{k \in U_2} \\ \vdots \\ \left[(x_{N_I} - \kappa_1)_+ \ (x_{N_I} - \kappa_2)_+ \ \cdots \ (x_{N_I} - \kappa_K)_+\right]_{k \in U_{N_I}} \end{bmatrix}_{N \times K}, \quad \boldsymbol{\beta}_2 = \begin{bmatrix} \beta_2 \\ \vdots \\ \beta_{1+K} \end{bmatrix},$$

$$Z_3 = \begin{bmatrix} [1 \ 0 \ \cdots \ 0]_{k \in U_1} \\ [0 \ 1 \ \cdots \ 0]_{k \in U_2} \\ \vdots \\ [0 \ 0 \ \cdots \ 1]_{k \in U_{N_I}} \end{bmatrix}_{N \times N_I} , \text{ and } \beta_3 = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_{N_I} \end{bmatrix} .$$

Now, write the spline fitting criterion as

$$\|Y - Z_1\beta_1 - Z_2\beta_2 - Z_3\beta_3\|^2 + \lambda_1 \|\beta_2\|^2 + \lambda_2 \|\beta_3\|^2 . \tag{5.5}$$

We can then use an argument similar to the one made in Chapter 1 to show that, up to a known multiplier, this is the same as the likelihood criterion obtained when finding the BLUP of $\beta_1$, $\beta_2$, and $\beta_3$ under the model assumption

$$Y|\beta_1, \beta_2, \beta_3 \sim N\left(Z_1\beta_1 + Z_2\beta_2 + Z_3\beta_3, \sigma^2 I\right), \quad \begin{pmatrix} \beta_2 \\ \beta_3 \end{pmatrix} \sim N\left(0, \mathbf{G}\right),$$

where

$$\mathbf{G} = \begin{bmatrix} \eta^2 I_{K \times K} & 0 \\ 0 & \tau^2 I_{N_I \times N_I} \end{bmatrix} .$$

Thus, the $(2 + K + N_I)$ vector $\beta = \left(\beta_1^T, \beta_2^T, \beta_3^T\right)^T$ that minimizes (5.5) is

$$B_U = \left(Z^T Z + D_\lambda\right)^{-1} Z^T Y \tag{5.6}$$

where

$$Z = [Z_1 \ Z_2 \ Z_3] = \left[\left[z_i^T\right]_{k \in U_i}\right]_{i \in U_I} \text{ and } D_\lambda = \text{diag}\left(0, 0, \underbrace{\lambda_1, \ldots, \lambda_1}_{K}, \underbrace{\lambda_2, \ldots, \lambda_2}_{N_I}\right) .$$

Let $m_i = N_i\mu_i = N_i z_i^T B_U$ be the model-predicted total for PSU $i$, assuming the entire population was observed. If these fitted values were known and if we were

able to obtain the cluster totals without error, then we would have a single stage sample and could estimate $t_y$ with a version of the *difference estimator* from Särndal et al. (1992, p. 221):

$$\hat{t}_{y,diff^*} = \sum_{i \in U_I} m_i + \sum_{i \in s_I} \frac{t_{yi} - m_i}{\pi_{Ii}}. \qquad (5.7)$$

However, we are not able to obtain the cluster totals without error. Due to cost or other considerations, a two-stage sampling design has been employed, and $t_{yi}$ is not observed for sampled clusters. In this case, we estimate $t_{yi}$ with its Horvitz-Thompson estimator

$$\hat{t}_{yi\pi} = \sum_{k \in s_i} \frac{y_k}{\pi_{k|i}},$$

which is design unbiased for $t_{yi}$ (see Horvitz and Thompson, 1952). What we will call the difference estimator is then obtained by (5.7), substituting $\hat{t}_{yi\pi}$ in for $t_{yi}$ in

$$\hat{t}_{y,diff} = \sum_{i \in U_I} m_i + \sum_{i \in s_I} \frac{\hat{t}_{yi\pi} - m_i}{\pi_{Ii}}. \qquad (5.8)$$

Both $\hat{t}_{y,diff^*}$ and $\hat{t}_{y,diff}$ are design unbiased for $t_y$. The design variance of (5.7) is

$$Var\left(\hat{t}_{y,diff^*}\right) = \sum_{i,j \in U_I} \sum \Delta_{Iij} \frac{t_{yi} - m_i}{\pi_{Ii}} \frac{t_{yj} - m_j}{\pi_{Ij}}, \qquad (5.9)$$

where $\Delta_{Iij} = \pi_{Iij} - \pi_{Ii}\pi_{Ij}$, and the design variance of (5.8) is

$$Var\left(\hat{t}_{y,diff}\right) = \sum_{i,j \in U_I} \sum \Delta_{Iij} \frac{t_{yi} - m_i}{\pi_{Ii}} \frac{t_{yj} - m_j}{\pi_{Ij}} + \sum_{i \in U_I} \frac{V_i}{\pi_{Ii}}, \qquad (5.10)$$

where

$$V_i = \sum_{k,l \in U_i} \sum \Delta_{kl|i} \frac{y_k}{\pi_{k|i}} \frac{y_l}{\pi_{l|i}},$$

and $\Delta_{kl|i} = \pi_{kl|i} - \pi_{k|i}\pi_{l|i}$. The extra term in (5.10) compared to (5.9) comes from the fact that we are estimating $t_y$.

Clearly, the estimator in (5.8) still cannot be calculated since we do not know $m_i$. Just as we replaced $t_{yi}$ with the sample-based estimator $\hat{t}_{yi\pi}$, so we will replace $m_i$ with its own sample-based estimator $\hat{m}_i$. To define $\hat{m}_i$, we follow Breidt et al. (2005); for a more detailed argument see Särndal et al. (1992, Chapter 8). First, define the diagonal matrix $\boldsymbol{W} = \mathrm{diag}_{k\in U}\{1/\pi_k\}$ and its sample submatrix $\boldsymbol{W}_s = \mathrm{diag}_{k\in s}\{1/\pi_k\}$. Also, let $\boldsymbol{Z}_s$ be the submatrix of $\boldsymbol{Z}$ consisting of those rows for which $k \in s$, and similarly for $\boldsymbol{Y}_s$ a subvector of $\boldsymbol{Y}$. For fixed $\lambda_1$ and $\lambda_2$ and under suitable regularity conditions, the weighted estimator

$$\hat{\boldsymbol{B}} = \left(\boldsymbol{Z}_s^T \boldsymbol{W}_s \boldsymbol{Z}_s + \boldsymbol{D}_\lambda\right)^{-1} \boldsymbol{Z}_s^T \boldsymbol{W}_s \boldsymbol{Y}_s = \boldsymbol{G}_\lambda \boldsymbol{Y}_s \qquad (5.11)$$

is a design-consistent estimator of $\boldsymbol{B}_U$. Finally, define $\hat{m}_i = N_i \hat{\mu}_i = N_i z_i^T \hat{\boldsymbol{B}}$. We can then define the model-assisted P-spline estimator as

$$\hat{t}_{y,spl} = \sum_{i\in U_I} \hat{m}_i + \sum_{i\in s_I} \frac{\hat{t}_{yi\pi} - \hat{m}_i}{\pi_{Ii}}. \qquad (5.12)$$

## 5.4  Properties of the estimator

Define the indicator function $I_{Ii} = 1$ if $i \in s_I$ and $I_{Ii} = 0$ otherwise, and the $n \times 1$ indicator vector $\mathbf{e}_i$ which is a vector of zeros except for a one at position $i$. Then, noting that $\boldsymbol{Y}_s = \sum_{k\in s} e_k y_k$, we can write (5.12) as

$$\hat{t}_{y,spl} = \sum_{i\in U_I} \hat{m}_i + \sum_{i\in s_I} \frac{\hat{t}_{yi\pi} - \hat{m}_i}{\pi_{Ii}}$$

$$
\begin{aligned}
&= \sum_{i \in U_I} N_i z_i^T \boldsymbol{G}_\lambda \boldsymbol{Y}_s + \sum_{i \in s_I} \frac{\sum_{k \in s_i} \frac{y_k}{\pi_{k|i}} - N_i z_i^T \boldsymbol{G}_\lambda \boldsymbol{Y}_s}{\pi_{Ii}} \\
&= \sum_{i \in U_I} N_i z_i^T \boldsymbol{G}_\lambda \boldsymbol{Y}_s + \sum_{i \in U_I} \sum_{k \in U_i} \frac{I_{Ii} I_{k|i}}{\pi_{Ii} \pi_{k|i}} y_k - \sum_{i \in U_I} \frac{I_{Ii}}{\pi_{Ii}} N_i z_i^T \boldsymbol{G}_\lambda \boldsymbol{Y}_s \\
&= \sum_{k \in s} \frac{y_k}{\pi_k} + \sum_{i \in U_I} N_i z_i^T \boldsymbol{G}_\lambda \left(1 - \frac{I_{Ii}}{\pi_{Ii}}\right) \sum_{k \in s} e_k y_k \\
&= \sum_{k \in s} \left[\frac{1}{\pi_k} + \sum_{i \in U_I} N_i z_i^T \boldsymbol{G}_\lambda e_k \left(1 - \frac{I_{Ii}}{\pi_{Ii}}\right)\right] y_k \\
&= \sum_{k \in s} w_k y_k
\end{aligned}
\tag{5.13}
$$

which shows that $\hat{t}_{y,spl}$ is a linear estimator, making it useful in a survey estimation context. If we let $\boldsymbol{t}_z = \sum_{i \in U_I} N_i \boldsymbol{z}_i$ and $\hat{\boldsymbol{t}}_{z,HT} = \sum_{i \in U_I} N_i \boldsymbol{z}_i / \pi_{Ii}$, then we can also write

$$
\begin{aligned}
\hat{t}_{y,spl} &= \sum_{i \in U_I} \hat{m}_i + \sum_{i \in s_I} \frac{\hat{t}_{yi\pi} - \hat{m}_i}{\pi_{Ii}} \\
&= \sum_{k \in s} \frac{y_k}{\pi_k} + \sum_{i \in U_I} N_i z_i^T \boldsymbol{G}_\lambda \left(1 - \frac{I_{Ii}}{\pi_{Ii}}\right) \boldsymbol{y} \\
&= \hat{t}_{y,HT} + \left(\boldsymbol{t}_z - \hat{\boldsymbol{t}}_{z,HT}\right)^T \hat{\boldsymbol{\beta}},
\end{aligned}
$$

which has the form of the traditional Horvitz-Thompson estimator plus a regression adjustment term.

### 5.4.1 Calibration

Särndal et al. (1992, Chapter 5) show that the ratio and linear regression estimators, in a single stage sampling design, have the property

$$
\hat{t} = \sum_{k \in U} \hat{m}_k.
\tag{5.14}
$$

Breidt et al. (2005) show this property also holds for the P-spline estimator in the single stage sampling design. Additionally, if an intercept is included in the model, they demonstrate that the P-spline estimator in the single stage sampling design is location and scale invariant, in the sense that

$$\sum_{k\in s} w_{k_s}\left(ay_k + b\right) = a\sum_{k\in s} w_{k_s}y_k + Nb \tag{5.15}$$

for any constants $a$ and $b$. This is a property shared by the general regression estimator.

In the two-stage sampling design, (5.14) and (5.15) do not necessarily hold. From (5.12), we see that (5.14) holds if $\sum_{i\in s_I} \hat{m}_i\pi_{Ii}^{-1} = \sum_{k\in s} y_k\pi_k^{-1}$. Further, the location invariance in (5.15) will hold provided that $\sum_{k\in s} w_{k_s} = N$. Proposition 1 and Proposition 2 address conditions in which (5.14) and (5.15) will hold, respectively, in the two-stage design.

**Proposition 1** *Consider the model-assisted P-spline estimator $\hat{t}_{y,spl}$ defined in (5.12). If*

$$\sum_{k\in s_i} \frac{1}{\pi_{k|i}} = N_i, \tag{5.16}$$

*then*

$$\sum_{i\in U_I} \frac{\hat{m}_i}{\pi_{Ii}} = \sum_{k\in s} \frac{y_k}{\pi_k}, \tag{5.17}$$

*and so $\hat{t}_{y,spl} = \sum_{i\in U_I} \hat{m}_i$.*

**Proof of Proposition 1:**

In the model specification (5.3), the existence of the intercept is important. This means that $\boldsymbol{e}_1^T \boldsymbol{z}_i = 1$ for all $i$ and $\boldsymbol{e}_1^T \mathbf{D}_\lambda = (0, 0, \ldots, 0)$,

$$
\begin{aligned}
\sum_{i \in s_I} \frac{\hat{m}_i}{\pi_{Ii}} &= \sum_{i \in s_I} \frac{N_i \boldsymbol{z}_i^T}{\pi_{Ii}} \left( \sum_{k \in s} \frac{\boldsymbol{z}_k \boldsymbol{z}_k^T}{\pi_k} + \mathbf{D}_\lambda \right)^{-1} \sum_{k \in s} \frac{\boldsymbol{z}_k y_k}{\pi_k} \\
&= \left( \sum_{i \in s_I} \frac{N_i \boldsymbol{e}_1^T \boldsymbol{z}_i \boldsymbol{z}_i^T}{\pi_{Ii}} + \boldsymbol{e}_1^T \mathbf{D}_\lambda \right) \left( \sum_{k \in s} \frac{\boldsymbol{z}_k \boldsymbol{z}_k^T}{\pi_k} + \mathbf{D}_\lambda \right)^{-1} \sum_{k \in s} \frac{\boldsymbol{z}_k y_k}{\pi_k} \\
&= \boldsymbol{e}_1^T \left( \sum_{i \in s_I} \frac{N_i \boldsymbol{z}_i \boldsymbol{z}_i^T}{\pi_{Ii}} + \mathbf{D}_\lambda \right) \left( \sum_{k \in s} \frac{\boldsymbol{z}_k \boldsymbol{z}_k^T}{\pi_k} + \mathbf{D}_\lambda \right)^{-1} \sum_{k \in s} \frac{\boldsymbol{z}_k y_k}{\pi_k} .
\end{aligned}
$$

(5.18)

If

$$
\left( \sum_{i \in s_I} \frac{N_i \boldsymbol{x}_i \boldsymbol{x}_i^T}{\pi_{Ii}} + \mathbf{D}_\lambda \right) \left( \sum_{k \in s} \frac{\boldsymbol{x}_k \boldsymbol{x}_k^T}{\pi_k} + \mathbf{D}_\lambda \right)^{-1} = \boldsymbol{I},
$$

(5.19)

then from (5.18) we get

$$
\begin{aligned}
\sum_{i \in s_I} \frac{\hat{m}_i}{\pi_{Ii}} &= \boldsymbol{e}_1^T \sum_{k \in s} \frac{\boldsymbol{x}_k y_k}{\pi_k} \\
&= \sum_{k \in s} \frac{y_k}{\pi_k} .
\end{aligned}
$$

To see when (5.19) holds, we write

$$
\begin{aligned}
\left( \sum_{k \in s} \frac{\boldsymbol{x}_k \boldsymbol{x}_k^T}{\pi_k} + \mathbf{D}_\lambda \right)^{-1} &= \left( \sum_{i \in s_I} \sum_{k \in s_i} \frac{\boldsymbol{x}_i \boldsymbol{x}_i^T}{\pi_{Ii} \pi_{k|i}} + \mathbf{D}_\lambda \right)^{-1} \\
&= \left( \sum_{i \in s_I} \frac{\boldsymbol{x}_i \boldsymbol{x}_i^T}{\pi_{Ii}} \sum_{k \in s_i} \frac{1}{\pi_{k|i}} + \mathbf{D}_\lambda \right)^{-1}
\end{aligned}
$$

Thus, if

$$
\sum_{k \in s_i} \frac{1}{\pi_{k|i}} = N_i,
$$

(5.20)

then (5.19) holds, and (5.14) follows.

One example of when the condition in (5.20) holds is in the case of a fixed size design within clusters with $\pi_{k|i} = n_i N_i^{-1}$. We now turn our attention to (5.15) with Proposition 2.

**Proposition 2** *Consider the model-assisted P-spline estimator $\hat{t}_{y,spl}$ defined in (5.12). If*

$$\sum_{k \in s_i} \frac{1}{\pi_{k|i}} = N_i, \tag{5.21}$$

*then the weights $\{w_{k_s}\}$ given in (5.13) satisfy*

$$\sum_{k \in s} w_{k_s} = N. \tag{5.22}$$

**Proof of Proposition 2:**

Again, using the fact that $e_1^T z_i = 1$ for all $i$ and $e_1^T D_\lambda = (0, 0, \ldots, 0)$, we have

$$
\begin{aligned}
\sum_{k \in s} w_k &= \sum_{k \in s} \left[ \frac{1}{\pi_k} + \sum_{i \in U_I} N_i z_i^T G e_k \left( 1 - \frac{I_{Ii}}{\pi_{Ii}} \right) \right] \\
&= \sum_{k \in s} \frac{1}{\pi_k} + \sum_{i \in U_I} \left( N_i z_i^T - \frac{I_{Ii} N_i z_i^T}{\pi_{Ii}} \right) \left( \sum_{k \in s} \frac{z_k z_k^T}{\pi_k} + D_\lambda \right)^{-1} \sum_{k \in s} \frac{z_k}{\pi_k} \\
&= \sum_{k \in s} \frac{1}{\pi_k} + \sum_{i \in U_I} \left( N_i z_i^T - \frac{I_{Ii} N_i z_i^T}{\pi_{Ii}} \right) \left( \sum_{k \in s} \frac{z_k z_k^T}{\pi_k} + D_\lambda \right)^{-1} \left( \sum_{k \in s} \frac{z_k z_k^T e_1}{\pi_k} + D_\lambda e_1 \right) \\
&= \sum_{k \in s} \frac{1}{\pi_k} + \sum_{i \in U_I} \left( N_i z_i^T - \frac{I_{Ii} N_i z_i^T}{\pi_{Ii}} \right) \left( \sum_{k \in s} \frac{z_k z_k^T}{\pi_k} + D_\lambda \right)^{-1} \left( \sum_{k \in s} \frac{z_k z_k^T}{\pi_k} + D_\lambda \right) e_1 \\
&= \sum_{k \in s} \frac{1}{\pi_k} + \sum_{i \in U_I} \left( N_i z_i^T - \frac{I_{Ii} N_i z_i^T}{\pi_{Ii}} \right) e_1 \\
&= \sum_{k \in s} \frac{1}{\pi_k} + \sum_{i \in U_I} N_i - \sum_{i \in s_I} \frac{N_i}{\pi_{Ii}} \\
&= \sum_{k \in s} \frac{1}{\pi_k} - \sum_{i \in s_I} \frac{N_i}{\pi_{Ii}} + N,
\end{aligned}
$$

so that (5.22) holds if and only if

$$\sum_{i \in s_I} \frac{N_i}{\pi_{Ii}} = \sum_{k \in s} \frac{1}{\pi_k}$$
$$= \sum_{i \in s_I} \sum_{k \in s_i} \frac{1}{\pi_{Ii}\pi_{k|i}}$$
$$= \sum_{i \in s_I} \frac{1}{\pi_{Ii}} \sum_{k \in s_i} \frac{1}{\pi_{k|i}}. \tag{5.23}$$

A sufficient condition for (5.23) is

$$\sum_{k \in s_i} \frac{1}{\pi_{k|i}} = N_i.$$

## 5.4.2 Asymptotic Results

Särndal et al. (1992) give an overview of the asymptotic framework for sampling from finite populations. More details and examples can be found in Isaki and Fuller (1982), Robinson and Särndal (1983), and Brewer (1979). In the two-stage sampling design, we consider an infinite sequence of sets of clusters $U_{I1} \subset U_{I2} \subset U_{I3} \subset \cdots$, where $U_{Iv}$ is a set of $N_{Iv}$ clusters of elements. Assume that $U_v = \bigcup_{i \in U_{Iv}} U_{vi}$, where cluster $U_{vi}$ contains $N_{vi} < \infty$ elements. We bound the number of elements in each cluster for every $v$ so that no one cluster dominates the population. It is assumed that $U_{vi} \subset U_{(v+1)i}$ for all $v$ and $i$. Then, as $v \to \infty$, $N_{Iv} \to \infty$. As a consequence, $N \to \infty$ although the number of elements in each cluster remains bounded.

Now, for each $U_{Iv}$, consider a first stage sampling design $p_{Iv}(\cdot)$ that assigns a certain probability $p_{Iv}(s_{Iv})$ to every possible sample of clusters $s_{Iv}$ from the set of clusters $U_{Iv}$. For simplicity, assume the first stage sample size, $n_{Iv}$, is fixed. Also, assume $n_{I1} < n_{I2} < n_{I3} < \cdots$. Thus, we also have $n_{Iv} \to \infty$ as $v \to \infty$. We

assume that the first stage sampling fraction $n_{Iv}/N_{Iv} \to \pi$ where $\pi \in (0,1)$. This is a necessary condition for assumption A5 below.

This framework, then, embeds the clusters $U_{Iv}$ and first stage sampling design $p_{Iv}(\cdot)$ into a sequence $\{U_{Iv}, p_{Iv}(\cdot)\}$ indexed by $v$. We could just as easily index such a sequence by $N_I$ and allow $N_I \to \infty$. This would mean that $n_I \to \infty$ in a predictable way if we make the assumption that $n_I/N_I \to \pi$ as $N_I \to \infty$. Thus, going forward, we drop the $v$ notation and consider the sequence $\{U_{I(N_I)}, p_{I(N_I)}(\cdot)\}$ indexed by $N_I$. The $o_p(\cdot)$ and $O_p(\cdot)$ notation below is with respect to the sequence of designs.

The following additional assumptions are also made:

**A1.** $B = \lim_{N_I \to \infty} B_U$ *exists, and* $\hat{B} - B_U = o_p(1)$. *Furthermore,*

$$E\left[\left(\hat{B}_i - B_{Ui}\right)^2 \left(\hat{B}_j - B_{Uj}\right)^2\right] = o(1) \text{ for } i,j = 1,\dots,p.$$

**A2.** *The limiting design covariance matrix of the normalized Horvitz-Thompson estimators,*

$$\Sigma = \begin{bmatrix} \Sigma_{yy} & \Sigma_{yz} \\ \Sigma_{yz}^T & \Sigma_{zz} \end{bmatrix} = \lim_{N_I \to \infty} \frac{n_I}{N_I^2} \begin{bmatrix} \sum\sum_{i,j \in U_I} \Delta_{Iij} \frac{t_i}{\pi_{Ii}} \frac{t_j}{\pi_{Ij}} + \sum_{i \in U_I} \frac{V_i}{\pi_{Ii}} & \sum\sum_{i,j \in U_I} \Delta_{Iij} \frac{t_{yi} z_j^T}{\pi_{Ii} \pi_{Ij}} \\ \sum\sum_{i,j \in U_I} \Delta_{Iij} \frac{z_i t_{yj}}{\pi_{Ii} \pi_{Ij}} & \sum\sum_{i,j \in U_I} \Delta_{Iij} \frac{z_i}{\pi_{Ii}} \frac{z_j^T}{\pi_{Ij}} \end{bmatrix}$$

*is positive definite, where*

$$V_i = \sum_{k \in U_i} \sum_{l \in U_i} \left(\pi_{kl|i} - \pi_{k|i} \pi_{l|i}\right) \frac{y_k y_l}{\pi_{k|i} \pi_{l|i}}$$

*and*

$$\Delta_{Iij} = \pi_{Iij} - \pi_{Ii} \pi_{Ij}.$$

**A3.** *The normalized Horvitz-Thompson estimators satisfy a central limit theorem:*

$$\frac{\sqrt{n_I}}{N_I} \begin{bmatrix} \sum_{i \in U_I} \left[\sum_{k \in U_i} \left(\frac{I_{Ii} I_{k|i}}{\pi_{Ii} \pi_{k|i}} - 1\right) y_k\right] \\ \sum_{i \in U_I} N_I z_i^T \left(\frac{I_{Ii}}{\pi_{Ii}} - 1\right) \end{bmatrix} \xrightarrow{dist} N(0, \Sigma)$$

as $N_I \to \infty$.

**A4.** *The estimated covariance matrix for the Horvitz-Thompson estimators is design*

*consistent in the following sense:*

$$\widehat{\Sigma} = \frac{n_I}{N_I^2} \left[ \begin{array}{cc} \sum \sum_{i,j \in s_I} \frac{\Delta_{Iij}}{\pi_{Iij}} \frac{\hat{t}_{yi\pi}}{\pi_{Ii}} \frac{\hat{t}_{yj\pi}}{\pi_{Ij}} + \sum_{i \in s_I} \frac{\hat{V}_i}{\pi_{Ii}} & \sum \sum_{i,j \in s_I} \frac{\Delta_{Iij}}{\pi_{Iij}} \frac{\hat{t}_{yi\pi} \mathbf{z}_j^T}{\pi_{Ii} \pi_{Ij}} \\ \sum \sum_{i,j \in s_I} \frac{\Delta_{Iij}}{\pi_{Iij}} \frac{\mathbf{z}_i \hat{t}_{yj\pi}}{\pi_{Ii} \pi_{Ij}} & \sum \sum_{i,j \in s_I} \frac{\Delta_{Iij}}{\pi_{Iij}} \frac{\mathbf{z}_i \, \mathbf{z}_j^T}{\pi_{Ii} \pi_{Ij}} \end{array} \right] \xrightarrow{p} \Sigma$$

as $N_I \to \infty$ , *where*

$$\hat{V}_i = \sum_{k \in s_i} \sum_{l \in s_i} \frac{\pi_{kl|i} - \pi_{k|i} \pi_{l|i}}{\pi_{kl|i}} \frac{y_k y_l}{\pi_{k|i} \pi_{l|i}}.$$

**A5.** *For all* $N_I$,

$$\min_{i \in U_I} \pi_{Ii} \geq \lambda > 0,$$

$$\min_{i,j \in U_I} \pi_{Iij} \geq \lambda^* > 0, \text{and}$$

$$\limsup_{N_I \to \infty} n_I \max_{i,j \in U_I : i \neq j} \left| \pi_{Iij} - \pi_{Ii} \pi_{Ij} \right| < \infty.$$

**A6.** *Additional assumptions involving higher-order inclusion probabilities:*

$$\limsup_{N_I \to \infty} n_I^2 \max_{(i_1,i_2,i_3,i_4) \in D_{4,N_I}} \left| E\left[ (I_{i_1} - \pi_{Ii_1})(I_{i_2} - \pi_{Ii_2})(I_{i_3} - \pi_{Ii_3})(I_{i_4} - \pi_{Ii_4}) \right] \right| = M_1 < \infty,$$

*where* $D_{t,N}$ *denotes the set of all distinct* $t$-*tuples* $(i_1, i_2, \ldots, i_t)$ *from* $U_I$,

$$\limsup_{N_I \to \infty} n_I \max_{(i_1,i_2,i_3) \in D_{3,N_I}} \left| E\left[ (I_{i_1} - \pi_{Ii_1})^2 (I_{i_2} - \pi_{Ii_2})(I_{i_3} - \pi_{Ii_3}) \right] \right| = M_2 < \infty,$$

*and*

$$\limsup_{N_I \to \infty} \max_{(i_1,i_2) \in D_{2,N_I}} \left| E\left[ (I_{i_1} - \pi_{Ii_1})^3 (I_{i_2} - \pi_{Ii_2}) \right] \right| = M_3 < \infty.$$

**A7.** *Additional assumptions involving element-wise products of the auxiliary vector:*

$$\lim_{N_I \to \infty} \frac{1}{N_I} \sum_{i=1}^{N_I} z_{ij}^4 < \infty \text{ for } j = 1, \ldots, p.$$

As Brewer (1979) notes, asymptotic analysis must allow the sample size, and therefore $\hat{t}_{y,spl}$, to go to infinity. Thus, we will consider the limit of $\hat{t}_{y,spl}/N$ rather than $\hat{t}_{y,spl}$.

The next series of results leads to a reasonable estimate for the variance of $\hat{t}_{y,spl}$ and an asymptotic distribution result using this estimate. Theorem 2 shows that a reasonable estimate for $Var\left(\hat{t}_{y,diff}\right)$ can be obtained by

$$\sum_{i\in s_I}\sum_{j\in s_I}\frac{\Delta_{Iij}}{\pi_{Iij}}\frac{\left(\hat{t}_{yi\pi}-\hat{m}_i\right)\left(\hat{t}_{yj\pi}-\hat{m}_j\right)}{\pi_{Ii}\pi_j}+\sum_{i\in s_I}\frac{\hat{V}_i}{\pi_{Ii}}. \tag{5.24}$$

**Theorem 2** *Under assumptions A1, A2, and A4,*

$$\sum_{i\in s_I}\sum_{j\in s_I}\frac{\Delta_{Iij}}{\pi_{Iij}}\frac{\left(\hat{t}_{yi\pi}-\hat{m}_i\right)\left(\hat{t}_{yj\pi}-\hat{m}_j\right)}{\pi_{Ii}\pi_j}+\sum_{i\in s_I}\frac{\hat{V}_i}{\pi_{Ii}}=Var\left(\hat{t}_{y,diff}\right)+o_p\left(\frac{N_I^2}{n_I}\right).$$

**Proof of Theorem 2:**

$$\sum_{i\in s_I}\sum_{j\in s_I}\frac{\Delta_{Iij}}{\pi_{Iij}}\frac{\left(\hat{t}_{yi\pi}-m_i\right)+\left(m_i-\hat{m}_i\right)}{\pi_{Ii}}\frac{\left(\hat{t}_{yj\pi}-m_j\right)+\left(m_j-\hat{m}_j\right)}{\pi_{Ij}}+\sum_{s_I}\frac{\hat{V}_i}{\pi_{Ii}}$$

$$=\sum_{i\in s_I}\sum_{j\in s_I}\frac{\Delta_{Iij}}{\pi_{Iij}}\frac{\left(\hat{t}_{yi\pi}-m_i\right)\left(\hat{t}_{yj\pi}-m_j\right)}{\pi_{Ii}\pi_{Ij}}+\sum_{s_I}\frac{\hat{V}_i}{\pi_{Ii}}$$

$$+\sum_{i\in s_I}\sum_{j\in s_I}\frac{\Delta_{Iij}}{\pi_{Iij}}\frac{\left(m_i-\hat{m}_i\right)\left(m_j-\hat{m}_j\right)}{\pi_{Ii}\pi_{Ij}}+2\sum_{i\in s_I}\sum_{j\in s_I}\frac{\Delta_{Iij}}{\pi_{Iij}}\frac{\left(\hat{t}_{yi\pi}-m_i\right)\left(m_j-\hat{m}_j\right)}{\pi_{Ii}\pi_{Ij}}$$

$$=\widehat{Var}\left(\hat{t}_{y,diff}\right)$$

$$+\sum_{i\in s_I}\sum_{j\in s_I}\frac{\Delta_{Iij}}{\pi_{Iij}}\frac{\left(m_i-\hat{m}_i\right)\left(m_j-\hat{m}_j\right)}{\pi_{Ii}\pi_{Ij}}+2\sum_{i\in s_I}\sum_{i\in s_I}\frac{\Delta_{Iij}}{\pi_{Iij}}\frac{\left(\hat{t}_{yi\pi}-m_i\right)\left(m_j-\hat{m}_j\right)}{\pi_{Ii}\pi_{Ij}}$$

$$=\widehat{Var}\left(\hat{t}_{y,diff}\right)+V_1+V_2 \tag{5.25}$$

By A4, $\widehat{Var}\left(\hat{t}_{y,diff}\right)=Var\left(\hat{t}_{y,diff}\right)+o_p\left(\frac{N_I^2}{n_I}\right)$. Also,

$$V_1=\left(\hat{B}-B_U\right)^T\sum_{i\in s_I}\sum_{i\in s_I}\frac{\Delta_{Iij}}{\pi_{Iij}}\frac{z_i\,z_j^T}{\pi_{Ii}\,\pi_{Ij}}\left(\hat{B}-B_U\right)$$

$$
\begin{aligned}
&= \; o_p\left(1\right) O_p\left(\frac{N_I^2}{n_I}\right) o_p\left(1\right) \\
&= \; o_p\left(\frac{N_I^2}{n_I}\right)
\end{aligned}
$$

and

$$
\begin{aligned}
V_2 \;&=\; 2 \sum_{i \in s_I} \sum_{j \in s_I} \frac{\Delta_{Iij}}{\pi_{Iij}} \frac{\bm{z}_i^T \bm{B}_U - \hat{t}_{yi\pi}}{\pi_{Ii}} \frac{\bm{z}_j^T}{\pi_{Ij}} \left(\hat{\bm{B}} - \bm{B}_U\right) \\
&=\; 2 \sum_{i \in s_I} \sum_{j \in s_I} \left[ \frac{\Delta_{Iij}}{\pi_{Iij}} \frac{\bm{B}_U \bm{z}_i \bm{z}_j^T}{\pi_{Ii}\pi_{Ij}} \left(\hat{\bm{B}} - \bm{B}_U\right) - \frac{\Delta_{Iij}}{\pi_{Iij}} \frac{\hat{t}_{yi\pi} \bm{z}_j^T}{\pi_{Ii}\pi_{Ij}} \left(\hat{\bm{B}} - \bm{B}_U\right) \right] \\
&=\; O_p\left(\frac{N_I^2}{n_I}\right) o_p\left(1\right) - O_p\left(\frac{N_I^2}{n_I}\right) o_p\left(1\right) \\
&=\; o_p\left(\frac{N_I^2}{n_I}\right)
\end{aligned}
$$

The result follows. †

We next show that, asymptotically, $Var\left(\hat{t}_{y,spl}\right)$ and $Var\left(\hat{t}_{y,diff}\right)$ are equivalent.

This will allow us to use the expression

$$
\sum_{i \in s_I} \sum_{j \in s_I} \frac{\Delta_{Iij}}{\pi_{Iij}} \frac{\left(\hat{t}_{yi\pi} - \hat{m}_i\right)\left(\hat{t}_{yj\pi} - \hat{m}_j\right)}{\pi_{Ii}\pi_{Ij}} + \sum_{i \in s_I} \frac{\hat{V}_i}{\pi_{Ii}}
$$

from Theorem 2 as an estimator for $Var\left(\hat{t}_{y,spl}\right)$. First, though, we need the following

two lemmas:

**Lemma 1** *Assumption A7 implies that the following results hold elementwise:*

$$
\lim_{N_I \to \infty} \frac{1}{N_I} \sum_{i=1}^{N_I} \bm{z}_i \bm{z}_i^T \bm{z}_i \bm{z}_i^T < \infty,
$$

$$
\lim_{N_I \to \infty} \frac{1}{N_I} \sum_{i=1}^{N_I} \bm{z}_i \bm{z}_i^T \bm{z}_i < \infty,
$$

$$
\lim_{N_I \to \infty} \frac{1}{N_I} \sum_{i=1}^{N_I} \bm{z}_i \bm{z}_i^T < \infty, \; and
$$

$$
\lim_{N_I \to \infty} \frac{1}{N_I} \sum_{i=1}^{N_I} \bm{z}_i < \infty.
$$

**Proof of Lemma 1:**

Write $\boldsymbol{z}_i = (z_{i1}, z_{i2} \ldots, z_{ip})^T$. Then,

$$
\boldsymbol{z}_i \boldsymbol{z}_i^T \boldsymbol{z}_i \boldsymbol{z}_i^T =
\begin{bmatrix}
z_{i1} \sum_{j=1}^p z_{ij}^2 z_{i1} & z_{i1} \sum_{j=1}^p z_{ij}^2 z_{i2} & \cdots & z_{i1} \sum_{j=1}^p z_{ij}^2 z_{ip} \\
z_{i2} \sum_{j=1}^p z_{ij}^2 z_{i1} & z_{i2} \sum_{j=1}^p z_{ij}^2 z_{i2} & \cdots & z_{i2} \sum_{j=1}^p z_{ij}^2 z_{ip} \\
\vdots & \vdots & \ddots & \vdots \\
z_{ip} \sum_{j=1}^p z_{ij}^2 z_{i1} & z_{ip} \sum_{j=1}^p z_{ij}^2 z_{i2} & \cdots & z_{ip} \sum_{j=1}^p z_{ij}^2 z_{ip}
\end{bmatrix},
$$

each element of which can be expressed as $z_{ia} \sum_{j=1}^p z_{ij}^2 z_{ib}$ for some $a, b \in \{1, \ldots, p\}$. Now since $\lim_{N_I \to \infty} \frac{1}{N_I} \sum_{i=1}^{N_I} z_{ij}^4 < \infty$ for $j = 1, \ldots, p$ by A7, we have by Cauchy-Schwartz that

$$
\begin{aligned}
\lim_{N_I \to \infty} \frac{1}{N_I} & \left( \sum_i z_{ia} \sum_j z_{ij}^2 z_{ib} \right) \\
&= \lim_{N_I \to \infty} \frac{1}{N_I} \sum_i \sum_j z_{ij}^2 z_{ia} z_{ib} \\
&\leq \lim_{N_I \to \infty} \left( \frac{\sum_i \sum_j z_{ij}^4}{N_I} \right)^{1/2} \left( \frac{\sum_i \sum_j z_{ia}^2 z_{ib}^2}{N_I} \right)^{1/2} \\
&\leq \lim_{N_I \to \infty} \left( \sum_j \frac{\sum_i z_{ij}^4}{N_I} \right)^{1/2} p^{1/2} \left( \frac{\sum_i z_{ia}^4}{N_I} \right)^{1/4} \left( \frac{\sum_i z_{ib}^4}{N_I} \right)^{1/4} \\
&< \infty.
\end{aligned}
$$

Thus, we have the following element-wise finite limits:

$$
\lim_{N_I \to \infty} \frac{1}{N_I} \sum_{i=1}^{N_I} \boldsymbol{z}_i \boldsymbol{z}_i^T \boldsymbol{z}_i \boldsymbol{z}_i^T < \infty.
$$

Similar Cauchy-Schwartz arguments show the remaining assertions in Lemma 1 hold.

†

**Lemma 2** *Under assumptions A5, A6, and Lemma 1,*

$$\lim_{N_I \to \infty} E \left\{ \left[ \frac{\sqrt{n_I}}{N_I} \sum_{i \in U_I} z_i \left( \frac{I_{Ii}}{\pi_{Ii}} - 1 \right) \right] \left[ \frac{\sqrt{n_I}}{N_I} \sum_{i \in U_I} z_i^T \left( \frac{I_{Ii}}{\pi_{Ii}} - 1 \right) \right] \right.$$

$$\left. \times \left[ \frac{\sqrt{n_I}}{N_I} \sum_{i \in U_I} z_i \left( \frac{I_{Ii}}{\pi_{Ii}} - 1 \right) \right] \left[ \frac{\sqrt{n_I}}{N_I} \sum_{i \in U_I} z_i^T \left( \frac{I_{Ii}}{\pi_{Ii}} - 1 \right) \right] \right\} < \infty, \quad (5.26)$$

*where the limits are element-wise.*

**Proof of Lemma 2:** We outline a sketch of the argument as straightforward bounding arguments and application of A5, A6, and Lemma 1 lead to the result.

Let $D_4$ be the set of all distinct 4-tuples of the set of indices $\{i, j, k, l\}$ and $D_3$ be the set of all distinct 3-tuples of the set of indices $\{i, j, k\}$. The expression in (5.26) consists of the sum of terms in the form of:

$$a_1 = \frac{n_I^2}{N_I^4} E \left[ \sum_{i,j,k,l \in U_I; i,j,k,l \in D_4} \frac{z_i z_j^T z_k z_l^T}{\pi_{Ii} \pi_{Ij} \pi_{Ik} \pi_{Il}} (I_{Ii} - \pi_{Ii}) (I_{Ij} - \pi_{Ij}) \right.$$

$$\left. \times (I_{Ik} - \pi_{Ik}) (I_{Il} - \pi_{Il}) \right]$$

$$a_2 = \frac{n_I^2}{N_I^4} E \left[ \sum_{i,j,k \in U_I; i,j,k \in D_3} \frac{z_i z_j^T z_k z_k^T}{\pi_{Ii} \pi_{Ij} \pi_{Ik}^2} (I_{Ii} - \pi_{Ii}) (I_{Ij} - \pi_{Ij}) (I_{Ik} - \pi_{Ik})^2 \right]$$

$$a_3 = \frac{n_I^2}{N_I^4} E \left[ \sum_{i,j \in U_I; i \neq j} \frac{z_i z_i^T z_j z_j^T}{\pi_{Ii}^2 \pi_{Ij}^2} (I_{Ii} - \pi_{Ii})^2 (I_{Ij} - \pi_{Ij})^2 \right]$$

$$a_4 = \frac{n_I^2}{N_I^4} E \left[ \sum_{i,j \in U_I; i \neq j} \frac{z_i z_j^T z_j z_j^T}{\pi_{Ii} \pi_{Ij}^3} (I_{Ii} - \pi_{Ii}) (I_{Ij} - \pi_{Ij})^3 \right]$$

$$a_5 = \frac{n_I^2}{N_I^4} E \left[ \sum_{i \in U_I} \frac{z_i z_i^T z_i z_i^T}{\pi_{Ii}^4} (I_{Ii} - \pi_{Ii})^4 \right] \quad (5.27)$$

Only a general form of each term is listed in (5.27). Now, by A5, A6, and Lemma 1,

$$a_1 \leq \frac{1}{N_I^4} \sum_{i,j,k,l \in U_I; i,j,k,l \in D_4} \frac{z_i z_j^T z_k z_l^T}{\pi_{Ii} \pi_{Ij} \pi_{Ik} \pi_{Il}} \times$$

$$n_I^2 \max \{ E \left[ (I_{Ii} - \pi_{Ii}) (I_{Ij} - \pi_{Ij}) (I_{Ik} - \pi_{Ik}) (I_{Il} - \pi_{Il}) \right] \}$$

$$
\begin{aligned}
&= M_1 \sum_{i \in U_I} \frac{\boldsymbol{z}_i}{N_I \pi_{Ii}} \sum_{j \in U_I; j \neq i} \frac{\boldsymbol{z}_j^T}{N_I \pi_{Ij}} \sum_{k \in U_I; k \neq i \neq j} \frac{\boldsymbol{z}_k}{N_I \pi_{Ik}} \sum_{l \in U_I; l \neq i \neq j \neq k} \frac{\boldsymbol{z}_l^T}{N_I \pi_{Il}} \\
&\leq \frac{M_1}{\lambda^4} \sum_{i \in U_I} \frac{\boldsymbol{z}_i}{N_I} \sum_{j \in U_I; j \neq i} \frac{\boldsymbol{z}_j^T}{N_I} \sum_{k \in U_I; k \neq i \neq j} \frac{\boldsymbol{z}_k}{N_I} \sum_{l \in U_I; l \neq i \neq j \neq k} \frac{\boldsymbol{z}_l^T}{N_I} \\
&= O(1),
\end{aligned}
$$

$$
\begin{aligned}
a_2 &\leq \frac{n_I}{N_I^4} \sum \sum \sum_{i,j,k \in U_I; i,j,k \in D_3} \frac{\boldsymbol{z}_i \boldsymbol{z}_j^T \boldsymbol{z}_k \boldsymbol{z}_k^T}{\pi_{Ii} \pi_{Ij} \pi_{Ik}^2} n_I \max \left\{ E\left[(I_{Ii} - \pi_{Ii})(I_{Ij} - \pi_{Ij})(I_{Ik} - \pi_{Ik})^2\right]\right\} \\
&\leq \frac{n_I M_2}{N_I \lambda^4} \sum_{i \in U_I} \frac{\boldsymbol{z}_i}{N_I} \sum_{j \in U_I; j \neq i} \frac{\boldsymbol{z}_j^T}{N_I} \sum_{k \in U_I; k \neq i \neq j} \frac{\boldsymbol{z}_k \boldsymbol{z}_k^T}{N_I} \\
&= O(1),
\end{aligned}
$$

$$
\begin{aligned}
a_3 &= \frac{n_I^2}{N_I^4} \sum \sum_{i,j \in U_I; i \neq j} \frac{\boldsymbol{z}_i \boldsymbol{z}_i^T \boldsymbol{z}_j \boldsymbol{z}_j^T}{\pi_{Ii}^2 \pi_{Ij}^2} E\left[(I_{Ii} - \pi_{Ii})^2 (I_{Ij} - \pi_{Ij})^2\right] \\
&\leq \frac{n_I^2}{N_I^2 \lambda^4} \sum_{i \in U_I} \frac{\boldsymbol{z}_i \boldsymbol{z}_i^T}{N_I} \sum_{j \in U_I; j \neq i} \frac{\boldsymbol{z}_j \boldsymbol{z}_j^T}{N_I} \\
&= O(1),
\end{aligned}
$$

$$
\begin{aligned}
a_4 &\leq \frac{n_I^2}{N_I^4} \sum \sum_{i,j \in U_I; i \neq j} \frac{\boldsymbol{z}_i \boldsymbol{z}_j^T \boldsymbol{z}_j \boldsymbol{z}_j^T}{\pi_{Ii} \pi_{Ij}^3} \max \left\{ E\left[(I_{Ii} - \pi_{Ii})(I_{Ij} - \pi_{Ij})^3\right]\right\} \\
&\leq \frac{n_I^2 M_3}{N_I^2 \lambda^4} \sum_{i \in U_I} \frac{\boldsymbol{z}_i}{N_I} \sum_{j \in U_I; j \neq i} \frac{\boldsymbol{z}_j^T \boldsymbol{z}_j \boldsymbol{z}_j^T}{N_I} \\
&= O(1), \text{ and}
\end{aligned}
$$

$$
\begin{aligned}
a_5 &= \frac{n_I^2}{N_I^4} \sum_{i \in U_I} \frac{\boldsymbol{z}_i \boldsymbol{z}_i^T \boldsymbol{z}_i \boldsymbol{z}_i^T}{\pi_{Ii}^4} E(I_{Ii} - \pi_{Ii})^4 \\
&\leq \frac{n_I^2}{N_I^2 \lambda^4} \sum_{i \in U_I} \frac{\boldsymbol{z}_i \boldsymbol{z}_i^T \boldsymbol{z}_i \boldsymbol{z}_i^T}{N_I^2} \\
&= O(1).
\end{aligned}
$$

Using similar bounding arguments on all terms of (5.26 gives the result.

**Theorem 3** *Under assumptions A1-A4,*

$$\lim_{N_I \to \infty} \frac{n_I}{N_I^2} Var\left(\hat{t}_{y,spl}\right) = \lim_{N_I \to \infty} \frac{n_I}{N_I^2} Var\left(\hat{t}_{y,diff}\right). \tag{5.28}$$

**Proof of Theorem 3:**

$$
\begin{aligned}
\frac{n_I}{N_I^2} Var\left(\hat{t}_{y,spl}\right) &= \frac{n_I}{N_I^2} Var\left[\sum_{i \in U_I} m_i + \sum_{i \in s_I} \frac{\hat{t}_{yi\pi} - m_i}{\pi_{Ii}} + \sum_{i \in U_I} (\hat{m}_i - m_i)\left(1 - \frac{I_{Ii}}{\pi_{Ii}}\right)\right] \\
&= \frac{n_I}{N_I^2} Var\left(\hat{t}_{y,diff}\right) + \frac{n_I}{N_I^2} Var\left[\sum_{i \in U_I} (\hat{m}_i - m_i)\left(1 - \frac{I_{Ii}}{\pi_{Ii}}\right)\right] \\
&\quad + \frac{n_I}{N_I^2} Cov\left[\hat{t}_{y,diff}, \sum_{i \in U_I} (\hat{m}_i - m_i)\left(1 - \frac{I_{Ii}}{\pi_{Ii}}\right)\right] \tag{5.29}
\end{aligned}
$$

The first term in (5.29) is

$$\frac{n_I}{N_I^2} Var\left(\hat{t}_{y,diff}\right) = \frac{n_I}{N_I^2}\left[\sum_{i \in U_I}\sum_{j \in U_I} \Delta_{Iij} \frac{t_{yi} - m_i}{\pi_{Ii}} \frac{t_{yj} - m_j}{\pi_{Ij}} + \sum_{i \in U_I} \frac{V_i}{\pi_{Ii}}\right],$$

which, by A2, converges to $\Sigma_{yy} + \boldsymbol{B}_U^T \Sigma_{zz} \boldsymbol{B}_U - 2\Sigma_{yz}\boldsymbol{B}_U$ as $N_I \to \infty$. The second term can be written as

$$
\begin{aligned}
\frac{n_I}{N_I^2} Var&\left[\sum_{i \in U_I} (\hat{m}_i - m_i)\left(1 - \frac{I_{Ii}}{\pi_{Ii}}\right)\right] \\
&= \frac{n_I}{N_I^2} Var\left[\sum_{i \in U_I} \left(\boldsymbol{z}_i^T \hat{\boldsymbol{B}} - \boldsymbol{z}_i^T \boldsymbol{B}_U\right)\left(1 - \frac{I_{Ii}}{\pi_{Ii}}\right)\right] \\
&= \frac{n_I}{N_I^2} Var\left[\sum_{i \in U_I} \boldsymbol{z}_i^T \left(1 - \frac{I_{Ii}}{\pi_{Ii}}\right)\left(\hat{\boldsymbol{B}} - \boldsymbol{B}_U\right)\right] \\
&\leq \frac{n_I}{N_I^2} E\left[\left(\sum_{i \in U_I} \boldsymbol{z}_i^T \left(1 - \frac{I_{Ii}}{\pi_{Ii}}\right)\left(\hat{\boldsymbol{B}} - \boldsymbol{B}_U\right)\right)^2\right] \\
&= \frac{n_I}{N_I^2} E\left[\left(\sum_{k=1}^{p}\sum_{i \in U_I} z_{ik}\left(1 - \frac{I_{Ii}}{\pi_{Ii}}\right)\left(\hat{B}_k - B_{Uk}\right)\right)^2\right]
\end{aligned}
$$

$$= \frac{n_I}{N_I^2} E \left[ \sum_{k=1}^{p} \sum_{l=1}^{p} \left( \sum_{i \in U_I} \sum_{j \in U_I} z_{ik} z_{jl} \left( 1 - \frac{I_{Ii}}{\pi_{Ii}} \right) \left( 1 - \frac{I_{Ij}}{\pi_{Ij}} \right) \right) \left( \hat{B}_k - B_{Uk} \right) \left( \hat{B}_l - B_{Ul} \right) \right]$$

$$= \frac{n_I}{N_I^2} \sum_{k=1}^{p} \sum_{l=1}^{p} E \left[ \left( \sum_{i \in U_I} \sum_{j \in U_I} z_{ik} z_{jl} \left( 1 - \frac{I_{Ii}}{\pi_{Ii}} \right) \left( 1 - \frac{I_{Ij}}{\pi_{Ij}} \right) \right) \left( \hat{B}_k - B_{Uk} \right) \left( \hat{B}_l - B_{Ul} \right) \right]$$

$$\leq \sum_{k=1}^{p} \sum_{l=1}^{p} \left\{ \frac{n_I^2}{N_I^4} E \left( \sum_{i \in U_I} \sum_{j \in U_I} z_{ik} z_{jl} \left( 1 - \frac{I_{Ii}}{\pi_{Ii}} \right) \left( 1 - \frac{I_{Ij}}{\pi_{Ij}} \right) \right)^2 \right\}^{1/2}$$

$$\times \left\{ E \left[ \left( \hat{B}_k - B_{Uk} \right) \left( \hat{B}_l - B_{Ul} \right) \right]^2 \right\}^{1/2}$$

$$= O(1) o(1) \tag{5.30}$$

by Lemma 2 and A1. Finally, because

$$\frac{n_I}{N_I^2} \left| Cov \left[ \hat{t}_{y,diff}, \sum_{i \in U_I} (\hat{m}_i - m_i) \left( 1 - \frac{I_{Ii}}{\pi_{Ii}} \right) \right] \right|$$

$$\leq \frac{n_I}{N_I^2} \sqrt{ Var \left[ \hat{t}_{y,diff} \right] Var \left[ \sum_{i \in U_I} (\hat{m}_i - m_i) \left( 1 - \frac{I_{Ii}}{\pi_{Ii}} \right) \right] }$$

$$\leq \sqrt{ \frac{n_I}{N_I^2} Var \left[ \hat{t}_{y,diff} \right] \frac{n_I}{N_I^2} Var \left[ \sum_{i \in U_I} (\hat{m}_i - m_i) \left( 1 - \frac{I_{Ii}}{\pi_{Ii}} \right) \right] }$$

$$\tag{5.31}$$

the third term in (5.29) equals $o(1)$, and the result is proven. †

The following theorem shows that $\hat{t}_{y,spl}$ is design consistent and provides its asymptotic distribution.

**Theorem 4** *Under assumptions A1-A3, the penalized spline estimator $\hat{t}_{y,spl}$ is design $\sqrt{n_I}$-consistent for $t_y$, in the sense that*

$$\frac{\hat{t}_{y,spl} - t_y}{N_I} = O_p \left( \frac{1}{\sqrt{n_I}} \right) \tag{5.32}$$

*and has the following distribution:*

$$\frac{\hat{t}_{y,spl} - t_y}{\sqrt{Var\left(\hat{t}_{y,diff}\right)}} \overset{dist}{\to} N\left(0,1\right) \tag{5.33}$$

**Proof of Theorem 4:**

$$
\begin{aligned}
\frac{\sqrt{n_I}}{N_I}\left(\hat{t}_{y,spl} - t_y\right) &= \frac{\sqrt{n_I}}{N_I}\left(\sum_{i\in U_I}\hat{m}_i + \sum_{i\in s_I}\frac{\hat{t}_{yi\pi} - \hat{m}_i}{\pi_{Ii}} - t_y\right) \\
&= \frac{\sqrt{n_I}}{N_I}\left[\sum_{i\in U_I}m_i + \sum_{i\in s_I}\frac{\hat{t}_{yi\pi} - m_i}{\pi_{Ii}} - t_y\right. \\
&\quad \left. + \sum_{i\in U_I}\hat{m}_i - \sum_{i\in U_I}m_i - \sum_{i\in s_I}\frac{\hat{m}_i}{\pi_{Ii}} + \sum_{i\in s_I}\frac{m_i}{\pi_{Ii}}\right] \\
&= \frac{\sqrt{n_I}}{N_I}\left[\left(\hat{t}_{y,diff} - ty\right) + \sum_{i\in U_I}N_iz_i^T\hat{B}\right. \\
&\quad \left. - \sum_{i\in U_I}N_iz_i^T\beta - \sum_{i\in U_I}I_{Ii}\frac{N_iz_i^T\hat{B}}{\pi_{Ii}} + \sum_{i\in U_I}I_{Ii}\frac{N_iz_i^T B_U}{\pi_{Ii}}\right] \\
&= \frac{\sqrt{n_I}}{N_I}\left[\left(\hat{t}_{y,diff} - ty\right) + \sum_{i\in U_I}N_iz_i^T\left(1 - \frac{I_i}{\pi_{Ii}}\right)\left(\hat{B} - B_U\right)\right] \\
&= \frac{\sqrt{n_I}}{N_I}\left(\hat{t}_{y,diff} - ty\right) + O_p\left(\frac{1}{\sqrt{n_I}}\right)o_p\left(1\right) \\
&= \frac{\sqrt{n_I}}{N_I}\left(\hat{t}_{y,diff} - ty\right) + o_p\left(\frac{1}{\sqrt{n_I}}\right) \tag{5.34}
\end{aligned}
$$

by A1 and A3. Focusing on the leading term in (5.34) we see that

$$
\begin{aligned}
\frac{\sqrt{n_I}}{N_I}\left(\hat{t}_{y,diff} - t_y\right) &= \frac{\sqrt{n_I}}{N_I}\left(\sum_{i\in U_I}m_i + \sum_{i\in s_I}\frac{\hat{t}_{yi\pi} - m_i}{\pi_{Ii}} - \sum_{k\in U}y_k\right) \\
&= \frac{\sqrt{n_I}}{N_I}\left(\sum_{i\in U_I}\frac{I_{Ii}}{\pi_{Ii}}\hat{t}_{yi\pi} - \sum_{k\in U}y_k + \sum_{i\in U_I}\frac{I_{Ii}}{\pi_{Ii}}m_i - \sum_{i\in U_I}m_i\right) \\
&= \frac{\sqrt{n_I}}{N_I}\left[\left(\sum_{k\in U}\frac{I_k}{\pi_k}y_k - \sum_{k\in U}y_k\right) + \left(\sum_{i\in U_I}z_i^T\left(\frac{I_{Ii}}{\pi_{Ii}} - 1\right)B_U\right)\right] \\
&= \frac{\sqrt{n_I}}{N_I}\left(\sum_{k\in U}\left(\frac{I_k}{\pi_k} - 1\right)y_k\right) + \frac{\sqrt{n_I}}{N_I}\left(\sum_{i\in U_I}z_i^T\left(\frac{I_{Ii}}{\pi_{Ii}} - 1\right)B_U\right).
\end{aligned}
$$

$$\tag{5.35}$$

Thus, by A2 and A3,

$$\frac{\sqrt{n_I}}{N_I}\left(\hat{t}_{y,diff} - t_y\right) \overset{dist}{\to} N\left(0, \Sigma_{yy} + \boldsymbol{B}_U^T \Sigma_{zz}\boldsymbol{B}_U - 2\Sigma_{yz}\boldsymbol{B}_U\right) = N\left(0, \tau^2\right),$$

which implies $\sqrt{n_I}\left(\hat{t}_{y,diff} - t_y\right)/N_I = O_p\left(1\right)$. From (5.34), we now get

$$\begin{aligned}
\frac{\hat{t}_{y,spl} - t_y}{N_I} &= O_p\left(\frac{1}{\sqrt{n_I}}\right) + o_p\left(\frac{1}{\sqrt{n_I}}\right) \\
&= O_p\left(\frac{1}{\sqrt{n_I}}\right),
\end{aligned}$$

which shows (5.32). To see that (5.33) holds, we note that

$$\begin{aligned}
Var\left(\frac{\sqrt{n_I}}{N_I}\left(\hat{t}_{y,diff} - t_y\right)\right) &= \frac{n_I}{N_I^2}Var\left(\hat{t}_{y,diff}\right) \\
&\to \tau^2 \text{ as } N_I \to \infty.
\end{aligned}$$

The result follows. †

## 5.5 Simulation Study

### 5.5.1 Simulation Design

A simulation study is conducted to compare the performance of the model-assisted estimator with P-spline to other choices. For the simulation, an artificial population is generated according to (5.1) with eight study variables whose mean functions are given by $m(x_{Ih}), h = 1, \ldots, N_I$, where $x_{Ih}$ is a known quantity associated with cluster $h$. Since auxiliary information is often proportional to cluster size, we use a design in which the first stage inclusion probabilities are proportional to the cluster size and let $x_{Ih} = \pi_{Ih}$. Eight different mean functions are generated similar to those from Breidt and Opsomer (2000):

Figure 5.1: Plots of the mean functions for the eight study variables used in the simulation.

linear: $m_1(\pi_{Ih}) = 10 + 5(\pi_{Ih} - 0.5),$

quadratic: $m_2(\pi_{Ih}) = 10 + 200(\pi_{Ih} - \overline{\pi_I})^2,$

bump: $m_3(\pi_{Ih}) = 10 + 5(\pi_{Ih} - 0.5) + \exp\left(-20000(\pi_{Ih} - 0.5)^2\right),$

jump: $m_4(\pi_{Ih}) = \left[10 + 5(\pi_{Ih} - 0.5) I_{\{\pi_{Ih} \leq 0.1\}}\right] + \left[11.25 + 5(\pi_{Ih} - 0.5) I_{\{\pi_{Ih} > 0.1\}}\right],$

exponential: $m_5(\pi_{Ih}) = 10 + \exp\left(-30(\pi_{Ih} - 0.05)\right),$

growth: $m_6(\pi_{Ih}) = 10 - \frac{\exp(-100(\pi_{Ih} - \overline{\pi_I}))}{1 + \exp(-100(\pi_{Ih} - \overline{\pi_I}))},$

cycle1: $m_7(\pi_{Ih}) = 10 + \sin\left(2\pi \frac{\pi_{Ih}}{\max(\pi_I)}\right),$

cycle4: $m_8(\pi_{Ih}) = 10 + \sin\left(8\pi \frac{\pi_{Ih}}{\max(\pi_I)}\right),$

Plots of the mean functions are shown in Figure 5.1.

We use Zheng and Little (2003) as a guideline for the structure of the population and for the sample design. The total number of Primary Sampling Units (PSU) is

fixed at 500; the size of the PSUs are integer values that can range from 50 to 400. The first stage sample is fixed at 48 PSUs, and the first order inclusion probabilities, $\pi_{Ii}, i = 1, \ldots, h$, are proportional to the size of the PSUs. The second stage inclusion probabilities, $\pi_{k|i}$, are inversely proportional to the first stage inclusion probabilities. This is an example of a *self-weighting design with a fixed first-stage sample*. In general, for such a design,

$$
\begin{aligned}
\pi_{Ii} &\propto N_i \\
&= c_1 N_i,
\end{aligned}
$$

$$
\begin{aligned}
\pi_{k|i} &\propto \frac{1}{\pi_{Ii}} \\
&= \frac{c_2}{\pi_{Ii}}, \text{ and}
\end{aligned}
$$

$$
\begin{aligned}
\pi_k &= \pi_{Ii} \pi_{k|i} \\
&= \pi_{Ii} \frac{c_2}{\pi_{Ii}} \\
&= c_2.
\end{aligned}
$$

From this we see that each element in the population has the same probability $\pi_k = c_2$ of being selected. Also, if we write $\pi_{k|i} = n_i / N_i$, then we see that $n_i = c_2 / c_1$, which is a constant. Thus, in a self-weighting design, the number of elements drawn from each cluster $i \in s_I$ is the same for every $i$.

## 5.5.2 The Estimators

Five different estimators through this simulation:

HT      Horvitz-Thompson

LIN     Model-assisted with simple linear model

SPL     Model-assisted with P-spline and no random effect for cluster

SPLRE   Model-assisted with P-spline and random effect for cluster

MBRE   Model-based with P-spline and random effect for cluster

The HT estimator is the $\pi$ estimator from Särndal et al. (1992, p. 137)

$$\hat{t}_{y\pi} = \sum_{s_I} \frac{\hat{t}_{yi\pi}}{\pi_{Ii}} = \sum_{s} \frac{y_k}{\pi_k}.$$

and is a generalization of the two-stage estimator given in Horvitz and Thompson

(1952). The three model-assisted estimators are computed as in (5.12); the only

difference is how $\hat{m}_i$ is computed. A summary of the model-assisted estimators

follows.

LIN

Working Model: $y_k = \beta_0 + \beta_1 \pi_{Ii} + \epsilon_k$

Estimate: $\hat{m}_i = N_i \left( \hat{\beta}_0 + \hat{\beta}_1 \pi_{Ii} \right)$

SPL

Working Model: $y_k = \beta_0 + \beta_1 \pi_{Ii} + \sum_{l=1}^{K} \beta_{1+l} \left( \pi_{ii} - \kappa_l \right)_+ + \epsilon_k$

Estimate: $\hat{m}_i = N_i \left( \hat{\beta}_0 + \hat{\beta}_1 \pi_{Ii} + \sum_{l=1}^{K} \hat{\beta}_{1+l} \left( \pi_{Ii} - \kappa_l \right)_+ \right)$

SPLRE

Working Model: $y_k = \beta_0 + \beta_1 x_i + \sum_{l=1}^{K} \beta_{1+l} \left( x_i - \kappa_l \right)_+ + \alpha_i + \epsilon_k$

Estimate: 
$$\hat{m}_i = \begin{cases} N_i \left( \hat{\beta}_0 + \hat{\beta}_1 \pi_{Ii} + \sum_{k=1}^{K} \hat{\beta}_{1+k} \left( \pi_{Ii} - \kappa_k \right)_+ + \hat{\alpha}_i \right) & \text{for } i \in s_I \\ N_i \left( \hat{\beta}_0 + \hat{\beta}_1 \pi_{Ii} + \sum_{k=1}^{K} \hat{\beta}_{1+k} \left( \pi_{Ii} - \kappa_k \right)_+ \right) & \text{otherwise.} \end{cases}$$

The MBRE estimator is taken from Zheng and Little (2003) and is computed as

$$\hat{t}_{yMB} = \sum_{i \in s_I} \left[ n_i \bar{y}_i + \left( N_i - n_i \right) \hat{\mu}_i \right] + \sum_{i \notin s_I} N_i \hat{\mu}_i, \tag{5.36}$$

where $\bar{y}_i$ is the sample mean of cluster $i \in s_I$ and

$$\hat{\mu}_i = \begin{cases} \hat{\beta}_0 + \hat{\beta}_1 \pi_{Ii} + \sum_{k=1}^{K} \hat{\beta}_{1+k} \left( \pi_{Ii} - \kappa_k \right)_+ + \hat{\alpha}_i & \text{for } i \in s_I \\ \hat{\beta}_0 + \hat{\beta}_1 \pi_{Ii} + \sum_{k=1}^{K} \hat{\beta}_{1+k} \left( \pi_{Ii} - \kappa_k \right)_+ & \text{otherwise.} \end{cases}$$

### 5.5.3 Knot selection and smoothing parameters

In Chapter 1, automatic knot selection and data driven smoothing parameters were discussed. We now consider what approach to take in the context of a large survey with many study variables. Mindful of the goal of obtaining one set of survey weights that can be applied to any study variable, it becomes clear that we need to set the number and location of knots and the smoothing parameters ahead of time. There are several ways this can be done.

If we had fit an unpenalized spline, then the number of knots and their location would have a dramatic effect on the fit. With too few knots, the fitted function may not be flexible enough to capture the complete trend in the data. With too many knots, the fit may be too rough; the large number of knots will allow the fitted function to change slope too often, which will cause it to pick up small fluctuations in the data that may not be of consequence. However, a penalized spline allows control over the effect of the knot terms through the smoothing parameter. Recall the penalty term from (1.16), $\lambda \beta^T D \beta$. Because of the matrix $D$, the penalty term penalizes only the coefficients of the knots, and if $\lambda$ is large, then the effect of the knots is diminished and the overall fit approaches the least-squares fit. Thus, one solution is to include a large number of knots in the model and allow the penalized least-squares criterion to determine the relative importance of each knot.

Ruppert et al. (2003) offer some guidelines in choosing the number of knots, noting that we want enough knots to resolve the underlying structure in the data but that too many knots will complicate computations. Through several demonstrations,

they show that there appears to be a threshold in the number of knots. Once this threshold is reached, adding more knots has a minimal effect on the fit, while a number of knots below this threshold leads to an over-smoothed fit. As a general rule, they suggest 4-5 unique observations between knots, with a maximum number of knots between 20 and 40. We will use the default choice for the number of knots $K$ given by Ruppert et al. (2003, p. 126)

$$K = \min \left( \frac{1}{4} \times \text{ number of unique } x_i, 35 \right). \tag{5.37}$$

The choice of the location of the knots will also affect the fit. More knots are needed in rougher areas of the scatter in order to capture the fluctuations with frequent slope changes, while smoother parts of the scatter will not need as many knots since there is less of a need for a change in slope. Ruppert et al. (2003, p. 126) suggest knot locations given by

$$\kappa_l = \left( \frac{l+1}{K+2} \right) th \text{ sample quantile of the unique } x_i, \; l = 1, \ldots, K. \tag{5.38}$$

Several data-driven methods for estimating the smoothing parameter are outlined in Chapter 1. Since the goal is to obtain one set of survey weights, we will not fit a smoothing parameter for each study variable. One option is to fix the degrees of freedom of the smooth and find the corresponding $\lambda$. A scatterplot smooth with $\nu$ degrees of freedom summarizes the data to about the same extent as a $(\nu - 1)$-degree polynomial, which provides some guideline in choosing degrees of freedom. A second reasonable option is, if one study variable is considered more important relative to the others, to use one of the data-driven methods for choosing $\lambda$ with

this study variable. This smoothing parameter estimate is then used for the other study variables. The P-spline with random effect for cluster in (5.3) requires two smoothing parameters, one for the knot terms and one for the random effect. We use a combination of the two. For the penalized spline with cluster random effect used in SPLRE and MBRE, we fit the model to the linear data using the function `lme` in `Splus`. Here the linear data is assumed to be the most important of the eight "study variables". The smoothing parameter associated with the random effect is estimated from this fit. The degrees of freedom due to the knot terms is then fixed and the corresponding degrees of freedom is computed as outlined in Chapter 1. A justification for combining these two methods for the two different smoothing parameters follows from Aerts et al. (2002).

### 5.5.4 Results

Population data were generated for each of the eight study variables with four different combinations of variance parameters $\sigma^2$ and $\tau^2$. One thousand samples were then drawn from each of these populations and the five estimators were computed for each sample. The estimators involving the penalized spline were computed with 4 degrees of freedom and with 10 degrees of freedom assigned to the knot terms. Mean square error (MSE) ratios are shown in Tables 5.1 and 5.2. The denominator in both tables is the MSE of the estimated sum of the response using the SPLRE method. Thus, ratios above 1 favor the use of SPLRE. It is clear that, in most cases, SPLRE does significantly better than HT, LIN, and SPL, while it often does as well

as MBRE. Recall that MBRE is a model-based method; it is, essentially, the result we would get by fitting a separate spline model to each study variable in each sample and ignoring the design properties of the simulation, instead relying solely on the model. Since the objective was to avoid having to rely on a model, we use MBRE as a comparison only and not as a viable option for our purposes. This re-fitting is often not an option in a real survey context, where a single set of weights is required.

## 5.6 Summary

In this chapter, we constructed a model-assisted survey estimator using a penalized spline. When the smoothing parameter is fixed, this estimator has the form of a linear estimator and also the form of the traditional Horvitz-Thompson estimator plus a regression adjustment term. Additionally, under reasonable conditions, the estimator is location and scale invariant. Asymptotically, $Var\left(\hat{t}_{y,spl}\right)$ can be estimated by (5.24). By fixing the degrees of freedom of the smooth obtained by the P-spline, we compute the weights $w_k$ in (5.13) once, regardless of the number of study variables. This yields a single set of weights that can be applied to all study variables. Fixing degrees of freedom of the smooth involves fixing the smoothing parameters and hence, the variance components; the variance components are not estimated for all study variables (or perhaps not for any study variable). Yet, the weights derived from $\hat{t}_{y,spl}$ will provide design consistent estimates for $t_y$.

The simulation study in Section 5.5 illustrates the usefulness of the estimator $\hat{t}_{y,spl}$. Data were generated for eight different study variables, all with different mean

Table 5.1: Mean square error ratios greater than one favor the model-assisted estimator with P-spline and cluster random effect. Results based on 1000 replicate random samples of size $n_I = 48$ and $n_i = 12$, $i \in s_I$. Degrees of freedom assigned to the knot terms is 4.

| Study variable | $\sigma^2$ | $\tau^2$ | HT | LIN | SPL | MBRE |
|---|---|---|---|---|---|---|
| linear | 0.01 | 0.01 | 35.65 | 1.65 | 1.95 | 1.01 |
| | 0.01 | 0.25 | 6.64 | 3.42 | 1.19 | 0.95 |
| | 0.25 | 0.01 | 1.40 | 1.25 | 1.00 | 1.04 |
| | 0.25 | 0.25 | 1.51 | 1.51 | 1.09 | 1.01 |
| quadratic | 0.01 | 0.01 | 29.92 | 43.73 | 3.12 | 0.40 |
| | 0.01 | 0.25 | 18.27 | 23.03 | 1.10 | 0.78 |
| | 0.25 | 0.01 | 1.10 | 1.38 | 1.00 | 0.97 |
| | 0.25 | 0.25 | 2.67 | 2.20 | 1.91 | 0.92 |
| bump | 0.01 | 0.01 | 21.19 | 3.13 | 0.99 | 0.97 |
| | 0.01 | 0.25 | 2.46 | 1.35 | 1.12 | 1.11 |
| | 0.25 | 0.01 | 1.43 | 1.05 | 1.00 | 0.98 |
| | 0.25 | 0.25 | 1.67 | 1.09 | 1.10 | 0.99 |
| jump | 0.01 | 0.01 | 10.87 | 4.13 | 1.27 | 0.79 |
| | 0.01 | 0.25 | 21.32 | 10.60 | 1.54 | 0.91 |
| | 0.25 | 0.01 | 3.56 | 1.97 | 1.00 | 0.93 |
| | 0.25 | 0.25 | 5.68 | 2.52 | 1.14 | 0.89 |
| exponential | 0.01 | 0.01 | 29.23 | 2.20 | 1.70 | 0.92 |
| | 0.01 | 0.25 | 13.63 | 1.89 | 1.26 | 0.91 |
| | 0.25 | 0.01 | 2.61 | 0.98 | 1.00 | 0.89 |
| | 0.25 | 0.25 | 3.76 | 0.94 | 1.07 | 0.97 |
| growth | 0.01 | 0.01 | 23.58 | 7.36 | 1.74 | 0.92 |
| | 0.01 | 0.25 | 16.51 | 8.19 | 1.35 | 1.00 |
| | 0.25 | 0.01 | 2.51 | 1.60 | 1.00 | 0.93 |
| | 0.25 | 0.25 | 3.47 | 1.68 | 1.06 | 0.99 |
| cycle1 | 0.01 | 0.01 | 12.89 | 5.31 | 2.03 | 0.79 |
| | 0.01 | 0.25 | 39.39 | 19.62 | 2.86 | 1.13 |
| | 0.25 | 0.01 | 1.97 | 1.06 | 1.01 | 0.89 |
| | 0.25 | 0.25 | 4.33 | 1.90 | 1.17 | 0.96 |
| cycle4 | 0.01 | 0.01 | 9.92 | 1.87 | 1.08 | 0.98 |
| | 0.01 | 0.25 | 1.56 | 1.61 | 1.76 | 0.91 |
| | 0.25 | 0.01 | 1.09 | 1.07 | 1.00 | 0.93 |
| | 0.25 | 0.25 | 1.12 | 1.21 | 1.02 | 1.08 |

functions. To obtain values for the smoothing parameters, we used a combination of data-driven methods and fixing the total number of degrees of freedom of the smooth while still only computing the $w_k$'s once. These weights were then applied to all of

Table 5.2: Mean square error ratios greater than one favor the model-assisted estimator with P-spline and cluster random effect. Results based on 1000 replicate random samples of size $n_I = 48$ and $n_i = 12$, $i \in s_I$. Degrees of freedom assigned to the knot terms is 10.

| Study variable | $\sigma^2$ | $\tau^2$ | HT | LIN | SPL | MBRE |
|---|---|---|---|---|---|---|
| linear | 0.01 | 0.01 | 15.94 | 1.14 | 1.16 | 0.97 |
| | 0.01 | 0.25 | 10.34 | 4.63 | 1.13 | 0.95 |
| | 0.25 | 0.01 | 1.69 | 1.29 | 1.34 | 0.99 |
| | 0.25 | 0.25 | 1.20 | 0.98 | 1.02 | 0.94 |
| quadratic | 0.01 | 0.01 | 28.46 | 9.20 | 1.07 | 0.91 |
| | 0.01 | 0.25 | 19.64 | 31.63 | 1.41 | 1.04 |
| | 0.25 | 0.01 | 3.61 | 2.48 | 1.06 | 0.97 |
| | 0.25 | 0.25 | 2.60 | 1.74 | 1.12 | 0.97 |
| bump | 0.01 | 0.01 | 7.27 | 2.68 | 1.73 | 0.72 |
| | 0.01 | 0.25 | 6.58 | 3.29 | 1.37 | 1.11 |
| | 0.25 | 0.01 | 1.34 | 1.11 | 1.07 | 1.02 |
| | 0.25 | 0.25 | 1.41 | 1.11 | 1.17 | 1.03 |
| jump | 0.01 | 0.01 | 10.94 | 10.38 | 2.54 | 0.87 |
| | 0.01 | 0.25 | 37.39 | 25.15 | 2.70 | 0.92 |
| | 0.25 | 0.01 | 4.55 | 2.48 | 1.12 | 0.95 |
| | 0.25 | 0.25 | 8.30 | 4.75 | 1.49 | 1.10 |
| exponential | 0.01 | 0.01 | 44.77 | 1.35 | 0.87 | 0.54 |
| | 0.01 | 0.25 | 39.47 | 1.96ˉ | 1.85 | 1.14 |
| | 0.25 | 0.01 | 2.72 | 0.94 | 1.30 | 1.07 |
| | 0.25 | 0.25 | 3.13 | 0.90 | 1.15 | 1.01 |
| growth | 0.01 | 0.01 | 12.49 | 4.20 | 1.28 | 0.93 |
| | 0.01 | 0.25 | 32.10 | 25.24 | 1.82 | 1.03 |
| | 0.25 | 0.01 | 2.80 | 1.68 | 1.20 | 1.04 |
| | 0.25 | 0.25 | 3.47 | 1.48 | 1.06 | 0.99 |
| cycle1 | 0.01 | 0.01 | 26.55 | 3.27 | 1.18 | 0.82 |
| | 0.01 | 0.25 | 32.01 | 18.80 | 1.37 | 1.05 |
| | 0.25 | 0.01 | 3.07 | 1.53 | 1.32 | 0.79 |
| | 0.25 | 0.25 | 2.97 | 2.11 | 1.23 | 1.03 |
| cycle4 | 0.01 | 0.01 | 32.96 | 3.52 | 1.17 | 0.87 |
| | 0.01 | 0.25 | 2.72 | 2.88 | 2.68 | 1.09 |
| | 0.25 | 0.01 | 1.02 | 1.10 | 1.04 | 0.91 |
| | 0.25 | 0.25 | 1.84 | 1.70 | 1.69 | 1.09 |

the study variables. The results using $\hat{t}_{y,spl}$ were compared to the Horvitz-Thompson estimator and to two other model-assisted estimators, one using a linear model and one using a P-spline with no random effect for cluster. Additionally, results from

$\hat{t}_{y,spl}$ were compared to the model-based estimator from Zheng and Little (2003). In the comparison to the model-based estimator, we fit a model to each study variable. As noted above, this was done for comparison. The results show that when the model-based estimator has a correctly specified model, it is usually better than $\hat{t}_{y,spl}$ in terms of MSE. However, $\hat{t}_{y,spl}$ is extremely competitive and is superior to the model-based estimator when the model is incorrectly specified.

Chapter 6

# DISCUSSION AND FUTURE RESEARCH

## 6.1 Summary of current research

We set forth to present and describe new methods that can be used in sample

surveys to estimate characteristics of domains of interest. These characteristics can

be specific quantities associated with an individual domain, as in estimating cluster

totals or mean ANC for a HUC, or properties of a collective set of domains, as in

the estimation of the distribution of change in ANC over time across all HUCs. This

dissertation introduced two methods for addressing estimation problems in sample

surveys: GCB estimators and a model-assisted estimator using penalized splines.

Chapter 2 introduced the GCB estimator and presented a numerical algorithm

for the computation of GCB estimates. The GCB estimates, like the CB estimates

introduced by Ghosh (1992), have good individual and ensemble properties. The

motivating example behind the GCB estimator was an ecological application. In

such a scenario, it is often of interest to obtain estimates on different scales, and

the ensemble properties of GCB estimates are examined under a transformation. We

demonstrate that the CB and GCB estimates of a transformation are not the CB and

GCB estimates, respectively, of the transformation. As an ensemble, the transformed

GCB estimates are better than the transformed Bayes estimates as, asymptotically, the edf of the transformed Bayes estimates is biased for the posterior mean of the edf of the transformed quantities of interest while, under reasonable conditions, the edf of the transformed GCB estimates appears not to be.

With computing power continuing to increase and with the availability of free (at the writing of this dissertation) software such as WinBUGS and R, Bayesian methods are readily accessible. Bayesian methods also can be an attractive approach to analysis with complex models since posterior means can be obtained for any parameter or function of parameters. In this context, the GCB estimator introduced in Chapter 2 should have a wide appeal in applications that require an understanding not only of individual domain characteristics, but also of ensemble characteristics.

Using the GCB estimator, we are able to extend the analysis of the Scotland lip cancer data from Stern and Cressie (1999) to a fully Bayesian context for a CAR model, producing results close to those of the earlier study. Having thus confirmed that GCB produces believable results for a CAR model, we apply the algorithm to a CAR model used on water quality data from the Mid-Atlantic Highlands. The nested structure of the watersheds in this region suggests that several levels of spatial correlation may exist. The GCB estimator is able to take advantage of this using a more complex model in which a CAR is placed on two levels of residuals. The flexibility of GCB is demonstrated in Chapter 4 where we reconsider the analysis of Opsomer et al. (2008) whose small area estimates for mean ANC in watersheds of the Northeast U.S. employ a penalized spline. As shown by Ruppert et al. (2003),

a P-spline can be formulated as a linear mixed model. Using the mixed-model specification we obtain GCB estimates that exhibit similar patterns to the results from Opsomer et al. (2008).

In Chapter 5, we looked at a second problem in survey estimation: obtaining a set of survey weights that can be applied to many study variables in a two-stage sampling design. The model-assisted estimator discussed in this chapter, $\hat{t}_{y,spl}$, was shown to have properties that make it useful in a survey context: it is a linear estimator, under reasonable conditions it is scale and location invariant and is calibrated to the model estimate, and it is design consistent. In a simulation study, $\hat{t}_{y,spl}$ was generally superior to other model-assisted estimators. Additionally, it was found to be competitive with a model-based estimator when the model was correctly specified and superior to the model-based estimator when the model was not correctly specified.

## 6.2  Future work

While the two methods discussed in this dissertation are fairly disparate, they are loosely tied together in their use of constrained or penalized methods. Future work includes an investigation of this connection in some detail. In particular, the Lagrangian function for the GCB estimator has the form of a sum of squares plus some "penalty" terms. This is similar to the form for maximum penalized likelihood in which we have a term to be minimized subject to a penalty for roughness. We want to examine what happens to the GCB estimates as the "penalty" coefficients

are allowed to vary and whether or not the GCB estimator has an interpretable mixed model representation similar to that of the penalized splines discussed in Chapters 1 and 5.

Currently, the GCB algorithm is a numerical procedure. It requires the estimation of a solution to a system of non-linear equations, and this system of equations can contain several unknown parameters. In order for the GCB algorithm to work well, the system of equations must have a stable solution and/or use a numerical algorithm that is able to converge correctly upon the solution. Our use of the R function `optim` to solve for Lagrangian parameters involves a minimization of squared error. In minimizing squared error, the algorithm converged more easily than when we solved the system of equations by minimizing absolute error. However, if the squared error function is flat in a neighborhood of the vertex, the numerical algorithm may have a more difficult time converging upon the correct minimum. Thus, options used in the numerical procedure should be investigated more thoroughly, including criteria for convergence and optimization methods. This is especially important when the Lagrangian is of high dimension. Use of a numerical algorithm to compute GCB estimates can be avoided if an analytical solution can be found.

In Chapter 2, we also noted that the GCB estimates of a function of the quantity of interest does not give the same results as applying the same function to the GCB estimates. Like CB estimators, GCB estimators constrain only low-order properties of the conditional distribution. In order for the ensemble properties of the GCB estimates to be preserved under a transformation, it is likely that at least some

constraint on the conditional covariance needs to be included. Depending upon the model, this may increase the dimensionality of the problem appreciably. However, using models for which the covariance structure can be specified relatively simply or with few parameters, such as an autoregressive process of order one, may make specification of a covariance constraint reasonable.

Chapter 5 focuses solely on one of the four cases describing the extent to which the auxiliary information is known and on what scale it is known, namely, Case A in which the auxiliary information is available for all clusters in the population. Recall that in Case D, auxiliary information is available for all clusters in the first-stage sample only. As such, it is difficult to make extensive use of the auxiliary information. Estimators for population quantities in this case must be design -based rather than model-based or model-assisted. Because of this, Case D is fairly uninteresting in the context of Chapter 5, and future work will focus on exploring Cases B and C. Of the two, Case B offers the most promise. In both Case B and C, regression modeling can be used to estimate quantities associated with elements. However, in Case C we can only obtain regression estimates of cluster-level quantities for clusters that are part of the first-stage sample. A design-based estimator must be used for quantities associated with the population. Case B is more conducive to the use of regression estimates since we have auxiliary information at the element-level from which cluster and population quantities can be estimated.

# Bibliography

Aerts, M., Claeskins, G., and Wand, M. (2002). Some theory for penalized spline generalized additive models. *Journal of Statistical Planning and Inference*, 103:455–470.

Banerjee, S., Carlin, B. P., and Gelfand, A. E. (2004). *Hierarchical Modeling and Analysis for Spatial Data*. Chapman & Hall, Boca Raton.

Besag, J. (1974). Spatial interaction and the statistical analysis of lattice systems. *Journal of the Royal Statistical Society. Series B*, 36:192–225.

Breidt, F. J., Claeskens, G., and Opsomer, J. D. (2005). Model-assisted estimation for complex surveys using penalized splines. *Biometrika*, 92:831–846.

Breidt, F. J. and Opsomer, J. D. (2000). Local polynomial regression estimators in survey sampling. *Annals of Statistics*, 28:1026–1053.

Brewer, K. R. W. (1979). A class of robust sampling designs for large-scale surveys. *Journal of the American Statistical Association*, 74:911–915.

Cassel, C. M., Särndal, C. E., and Wretman, J. (1977). *Foundations of Inference in Survey Sampling*. Wiley, New York.

Clayton, D. G. and Kaldor, J. (1987). Empirical Bayes estimats of age-standardized relative risks for use in disease mapping. *Biometrics*, 43:671–681.

Cordy, C. B. and Thomas, D. R. (1997). Deconvolution of a distribution function. *Journal of the American Statistical Association*, 92:1459–1465.

Cressie, N. and Chan, N. H. (1989). Spatial modeling of regional variables. *Journal of the American Statistical Association*, 84:393–401.

Cressie, N. A. C. (1993). *Statistics for Spatial Data*. John Wiley & Sons, New York.

Devine, O. J. and Louis, T. A. (1994). A constrained empirical Bayes estimator for incidence rates in areas with small populations. *Statistics in Medicine*, 13:1119–1133.

Devine, O. J., Louis, T. A., and Halloran, M. E. (1994a). Empirical Bayes estimators for spatially correlated incidence rates. *Evironmentrics*, 5:381–398.

Devine, O. J., Louis, T. A., and Halloran, M. E. (1994b). Empirical Bayes methods for stabilizing incidence rates before mapping. *Epidemiology*, 5:622–630.

Eilers, P. H. C. and Marx, B. D. (1996). Flexible smoothing with b-splines and penalties (with discussion). *Statistical Science*, 11:89–121.

Fay, R. E. and Herriot, R. A. (1979). Estimation of income from small places: An application of James-Stein procedures to census data. *Journal of the American Statistical Association*, 74:269–277.

Gelman, A., Carlin, J. B., Stern, H. S., and Rubin, D. B. (2004). *Bayesian Data Analysis*. Chapman & Hall, Boca Raton, second edition.

Gelman, A. and Price, P. N. (1999). All maps of parameter estimates are misleading. *Statistics in Medicine*, 18:3221–3234.

Gelman, A. and Rubin, D. B. (1992). Inference from iterative simulation using multiple sequences (with discussion). *Statistical Science*, 17:457–511.

Ghosh, M. (1992). Constrained Bayes estimation with applications. *Journal of the American Statistical Association*, 87:533–540.

Green, P. J. and Silverman, B. W. (1994). *Nonparamentric regression and generalized linear models*. Chapman & Hall, Boca Raton.

Hardle, W., Kerkyacharian, G., Picard, D., and Tsybakov, A. (1998). *Wavelets, Approximation, and Statistical Applications (Lecture Notes in Statistics, vol. 29)*. Springer-Verlag, New York.

Hocking, R. (1996). *Methods and Applications of Linear Models*. John Wiley & Sons, New York.

Horvitz, D. G. and Thompson, D. J. (1952). A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association*, 47:663–685.

Isaki, C. T. and Fuller, W. A. (1982). Survey design under the regression superpopulation model. *Journal of the American Statistical Association*, 77:89–96.

Kleinschmidt, I., Omumbo, J., Briet, O., van de Giesen, N., Sogoba, N., Mensah, N. K., Windmeijer, P., Moussa, M., and Teuscher, T. (2001). An empirical malaria distribution map for West Africa. *Tropical Medicine & International Health*, 6:779–786.

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.

Lunn, D. J., Thomas, A., Best, N., and Spiegelhalter, D. (2000). Winbugs – a Bayesian modelling framework: concepts, structure, and extensibility. statistics and computing. *Statistics and Computing*, 10:325–337.

Manton, K. G., Woodbury, M. A., Stallard, E., Riggan, W. B., Creason, J. P., and Pellom, A. C. (1989). Empirical Bayes procedures for stabilizing maps of U.S. cancer mortality rates. *Journal of the American Statistical Association*, 84:637–650.

McCulloch, C. E. and Searle, S. R. (2001). *Generalized, Linear and Mixed Models*. Wiley, New York.

Nychka, D. and Cummins, D. (1996). Commment on paper by Eilers and Marx. *Statistical Science*, 11:104–105.

Ogden, R. T. (1996). *Essential Wavelets for Statistical Applications and Data Analysis*. Birkhauser, Boston.

Opsomer, J. D., Claeskens, G., Ranalli, M. G., Kauermann, G., and Breidt, F. J. (2008). Nonparametric small area estimation using penalized spline regression. *Journal of the Royal Statistical Association, Series B*, 70:265–286.

Rao, J. N. K. (2003). *Small Area Estimation*. John Wiley & Sons, Hoboken.

Robinson, P. M. and Särndal, C. E. (1983). Asymptotic properties of the generalized regression estimator in probability sampling. *Sankhya B*, 45:240–248.

Ruppert, D., Wand, M. P., and Carroll, R. J. (2003). *Semiparametric Regression*. Cambridge University Press, Cambridge, UK.

Särndal, C. E., Swensson, B., and Wretman, J. (1992). *Model Assisted Survey Sampling*. Springer-Verlag, New York.

Seaber, P. R., Kapinos, F. P., and Knapp, G. L. (1987). Hydrologic unit maps. U.S. Geological Survey Water-Supply Paper 2294, United States Department of the Interior, Geological Survey, Denver, CO.

Stern, H. and Cressie, N. (1999). *Disease Mapping and Risk Assessment for Public Health*, chapter 5, pages 63–84. John Wiley & Sons, Chichester.

Stoddard, J. L., Kahl, J. S., Deviney, F. A., DeWalle, D. R., Driscoll, C. T., Herlihy, A. T., Kellogg, J. H., Murdoch, P. S., Webb, J. R., and Webster, K. E. (2002). Response of surface water chemistry to the Clean Air Act Amendments of 1990. Technical report, Environmental Protection Agency, Corvallis, OR.

Zheng, H. and Little, R. (2003). Inference for the population total from probability-proportional-to-size samples based on predictions from a penalized spline non-parametric model. Working paper 9, Department of Biostatistics, University of Michigan, Ann Arbor, MI.