

An Approximate Bayesian Approach to Regression Estimation with Many Auxiliary Variables

Shonosuke Sugasawa¹ Jae Kwang Kim²

June 12, 2019

Abstract

Model-assisted estimation with complex survey data is an important practical problem in survey sampling. When there are many auxiliary variables, selecting significant variables associated with the study variable would be necessary to achieve efficient estimation of population parameters of interest. In this paper, we formulate a regularized regression estimator in the framework of Bayesian inference using the penalty function as the shrinkage prior for model selection. The proposed Bayesian approach enables us to get not only efficient point estimates but also reasonable credible intervals for population means. Results from two limited simulation studies are presented to facilitate comparison with existing frequentist methods.

Keywords: Generalized regression estimation, Regularization, Shrinkage prior, Survey Sampling

1 Introduction

Probability sampling is a scientific tool for obtaining a representative sample from the target population. In order to estimate a finite population total from a target population, the Horvitz-Thompson estimator obtained from a probability sample is often used, which satisfies consistency and the resulting inference is justified from the randomization perspective (Horvitz and Thompson, 1952). However, the Horvitz-Thompson estimator uses the first-order inclusion probability only and

¹ Center for Spatial Information Science, The University of Tokyo

² Department of Statistics, Iowa State University

does not fully incorporate all available information from the finite population. To improve efficiency, regression estimation is often used to incorporate auxiliary information in survey sampling. Deville and Saˆrndal (1992), Fuller (2002), Kim and Park (2010), and Breidt and Opsomer (2017) present comprehensive overviews of such variants of regression estimation in survey sampling.

The regression estimation approaches in survey sampling assume a model for the finite population, i.e., the superpopulation model, as

$$y_i = x_i^t \beta + e_i, \quad (1)$$

where $E(e_i) = 0$ and $\text{Var}(e_i) = \sigma^2$. The superpopulation model does not necessarily hold in the sample as the sampling design can be informative in the sense of Pfeffermann and Sverchkov (1999). Under the regression superpopulation model in (1), Isaki and Fuller (1982) show that the asymptotic variance of the regression estimator achieves the lower bound of Godambe and Joshi (1965). Thus, the regression estimator is asymptotically efficient in the sense of achieving the minimum variance under the joint distribution of the sampling design and the superpopulation model in (1).

However, the above optimality of the regression estimator is untenable if the dimension of the auxiliary variables x is large. When there are many auxiliary variables, the asymptotic bias of the regression estimator using all the auxiliary variables is no longer negligible and the resulting inference can be problematic.

Simply put, including irrelevant auxiliary variables can introduce substantial variability in point estimation, but the uncertainty can fail to be fully accounted for by the standard linearization variance estimation, resulting in misleading inference.

To overcome the problem, Saˆrndal and Lundstroˆm (2005) select a subset of the auxiliary variables for regression estimation. The classical model selection

approach is based on a step-wise method. However, the step-wise methods will not necessarily produce the best model if there are redundant predictors. Another approach is to employ regularized estimation of regression coefficients. For example, McConville et al. (2017) propose a regularized regression estimation approach based on the LASSO penalty of Tibshirani (1996). However, there are two main problems in the regularization approach. First, the choice of the regularization parameter is somewhat unclear. Second, after model selection, the frequentist inference is notoriously difficult to make.

In this paper, we propose a unified Bayesian framework to handle regularized regression estimation. We first present a Bayesian approach for regression estimation when $p = \dim(x)$ is fixed, using the approximate Bayesian approach considered in Wang et al. (2018). The proposed Bayesian method fully captures the uncertainty in parameter estimation for the regression estimator and has better coverage properties. Second, the proposed Bayesian method solves the problem of large p in regularized regression estimation.

The penalty function for regularization is incorporated into the prior distribution and the uncertainty associated with model selection and parameter estimation is fully captured in the Bayesian machinery. Furthermore, the penalty parameter λ can be optimized by having its own prior distribution. The proposed method provides a unified approach to Bayesian inference with sparse regression estimation. It is a calibrated Bayesian (Little, 2012) in the sense that it is asymptotically equivalent to the frequentist design-based approach.

The paper is organized as follows. In Section 2, the basic setup is introduced. In Section 3, the approximate Bayesian inference using regression estimation is proposed under a fixed p . In Section 4, the proposed method is extended to handle sparse regression estimation using shrinkage prior distributions. In Section 5, the proposed method is extended to non-linear regression models. In Section 6, results from two limited simulation studies are presented. The proposed method is

applied to the real data example in Section 7. Some concluding remarks are made in Section

8.

2 Basic setup

Consider a finite population of a known size N . Associated with unit i in the finite population, we consider measurement (x_i^t, y_i) where x_i is the vector of auxiliary variables with dimension p and y_i is the study variable of interest. We are interested in estimating the finite population mean $\bar{Y} = N^{-1} \sum_{i=1}^N y_i$ from a sample selected by a probability sampling design. Let A be the index set of the sample and we observe $\{x_i, y_i\}_{i \in A}$ from the sample. The Horvitz-Thompson estimator $\hat{Y}_{HT} = N^{-1} \sum_{i \in A} \pi_i^{-1} y_i$ is design unbiased but it is not necessarily efficient.

If the finite population mean $\bar{X} = N^{-1} \sum_{i=1}^N x_i$ is known, then we can improve the efficiency of \hat{Y}_{HT} by using the following regression estimator:

$$\hat{Y}_{\text{reg}} = \bar{X}^t \hat{\beta} + \frac{1}{N} \sum_{i \in A} \frac{1}{\pi_i} (y_i - x_i^t \hat{\beta}) \quad (2)$$

where π_i is the first-order inclusion probability of unit i , and $\hat{\beta}$ is an estimator of β in (1). Typically, we use $\hat{\beta}$ obtained by minimizing the weighted quadratic loss

$$Q(\beta) = \sum_{i \in A} \pi_i^{-1} (y_i - x_i^t \beta)^2, \quad (3)$$

motivated from model (1).

To derive the asymptotic properties of \hat{Y}_{reg} , we may use

$$\begin{aligned} \hat{Y}_{\text{reg}} - \bar{Y} &= \hat{Y}_{HT} - \bar{Y} + \left(\bar{X} - \hat{X}_{HT} \right)^t \hat{\beta} \\ &= \hat{Y}_{HT} - \bar{Y} + \left(\bar{X} - \hat{X}_{HT} \right)^t \beta_* + R_n \end{aligned} \quad (4)$$

where $\mathbf{X}_{\text{HT}}^{-1} = N^{-1} \sum_{i \in A} \pi_i^{-1} \mathbf{x}_i$ and

$$R_n = \left(\bar{\mathbf{X}} - \hat{\mathbf{X}}_{\text{HT}} \right)^t \left(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_* \right)$$

for any $\boldsymbol{\beta}_*$. If we choose $\boldsymbol{\beta}_* = \text{plim}_{n \rightarrow \infty} \hat{\boldsymbol{\beta}}$ and the dimension p is fixed in the asymptotic setup, then we can obtain $R_n = O_p(n^{-1})$ and safely use the main terms of (4) to describe the asymptotic behavior of \hat{Y}_{reg}^{-} . To emphasize its dependence on $\hat{\boldsymbol{\beta}}$ in the regression estimator, we can write $\hat{Y}_{\text{reg}}^{-} = \hat{Y}_{\text{reg}}^{-}(\hat{\boldsymbol{\beta}})$. Roughly speaking, we can obtain

$$\sqrt{n} \left(\hat{Y}_{\text{reg}}^{-}(\hat{\boldsymbol{\beta}}) - \hat{Y}_{\text{reg}}^{-}(\boldsymbol{\beta}_*) \right) = O_p(n^{-1/2} p) \quad (5)$$

and, if $p = o(n^{1/2})$ then we can safely ignore the effect of estimating $\boldsymbol{\beta}_*$ in the regression estimator. See Appendix A for a sketch proof of (5).

If, on the other hand, the dimension p is large, then we cannot ignore the effect of estimating $\boldsymbol{\beta}_*$. In this case, we can consider using some variable selection idea to reduce the dimension of \mathbf{X} . For variable selection, we may employ techniques of regularized estimation of regression coefficients. The regularization method can be described as finding

$$\hat{\boldsymbol{\beta}}^{(R)} = \underset{\boldsymbol{\beta}}{\text{argmin}} \{ Q(\boldsymbol{\beta}) + p_\lambda(\boldsymbol{\beta}) \}, \quad (6)$$

where $Q(\boldsymbol{\beta})$ is defined in (3) and $p_\lambda(\boldsymbol{\beta})$ is a penalty function with parameter λ . Some popular penalty functions are presented in Table 1. Once the solution to (6) is obtained, then the regularized regression estimator is given by

$$\hat{Y}_{\text{reg}}(\hat{\boldsymbol{\beta}}^{(R)}) = \bar{\mathbf{X}}^t \hat{\boldsymbol{\beta}}^{(R)} + \frac{1}{N} \sum_{i \in A} \frac{1}{\pi_i} \left(y_i - \mathbf{x}_i^t \hat{\boldsymbol{\beta}}^{(R)} \right). \quad (7)$$

Table 1: Popular penalized regression methods

Method	Reference	Penalty function
Ridge	Hoerl and Kennard (1970)	
LASSO	Tibshirani (1996)	
Adaptive LASSO	Zou (2006)	$p_\lambda(\boldsymbol{\beta}) = \lambda \sum_{j=1}^p \beta_j^2$ $p_\lambda(\boldsymbol{\beta}) = \lambda \sum_{j=1}^p \beta_j $ $p_\lambda(\boldsymbol{\beta}) = \lambda \sum_{j=1}^p \left(\beta_j / \left \hat{\beta}_j \right \right)$
Elastic Net	Zou and Hastie (2005)	$p_\lambda(\boldsymbol{\beta}) = \lambda_1 \sum_{j=1}^p \beta_j + \lambda_2 \sum_{j=1}^p \beta_j^2$

Statistical inference with the regularized regression estimator in (7) is not fully investigated in the literature. For example, Chen et al. (2018) consider the regularized regression estimator using adaptive LASSO of Zou (2006), but they assume the sampling design is non-informative and the uncertainty in model selection is not fully incorporated in their inference. Generally speaking, making frequentist inference after model selection is difficult. The approximated Bayesian method we propose in this paper will capture the full uncertainty in the Bayesian framework.

3 Approximate Bayesian survey regression estimation

Developing a design-based Bayesian inference under complex sampling is a challenging problem in statistics. Wang et al. (2018) propose the so-called approximate Bayesian method for design-based inference using asymptotic normality of a designconsistent estimator. Specifically, for a given parameter θ with a prior distribution $\pi(\theta)$, if one can find a design-consistent estimator $\hat{\theta}$ of θ , then the approximate posterior distribution of θ is given by

$$p(\theta | \hat{\theta}) = \frac{f(\hat{\theta} | \theta)\pi(\theta)}{\int f(\hat{\theta} | \theta)\pi(\theta)d\theta}, \quad (8)$$

where $f(\hat{\theta} \mid \theta)$ is the sampling distribution of $\hat{\theta}$, which is often approximated by a normal distribution.

Drawing on this idea, one can develop an approximate Bayesian approach to capture the full uncertainty in the regression estimator. Let

$$\hat{\beta} = \left(\sum_{i \in A} \pi_i^{-1} \mathbf{x}_i \mathbf{x}_i^t \right)^{-1} \sum_{i \in A} \pi_i^{-1} \mathbf{x}_i y_i$$

be the design-consistent estimator of β and \hat{V}_β be the corresponding asymptotic variance-covariance matrix of $\hat{\beta}$, given by

$$\hat{V}_\beta = \left(\sum_{i \in A} \pi_i^{-1} \mathbf{x}_i \mathbf{x}_i^t \right)^{-1} \left(\sum_{i \in A} \sum_{j \in A} \frac{\Delta_{ij}}{\pi_{ij}} \frac{\hat{e}_i \mathbf{x}_i}{\pi_i} \frac{\hat{e}_j \mathbf{x}_j^t}{\pi_j} \right) \left(\sum_{i \in A} \pi_i^{-1} \mathbf{x}_i \mathbf{x}_i^t \right)^{-1}, \quad (9)$$

where $\hat{e}_i = y_i - \mathbf{x}_i^t \hat{\beta}$, $\Delta_{ij} = \pi_{ij} - \pi_i \pi_j$ and π_{ij} is the joint inclusion probability of unit i and j . Under some regularity conditions, as discussed in Chapter 2 of Fuller (2009), we can establish

$$\hat{V}_\beta^{-1/2} (\hat{\beta} - \beta) \mid \beta \xrightarrow{\mathcal{L}} N(0, I_p) \quad (10)$$

as $n \rightarrow \infty$, where

$$\beta = \left(\sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^t \right)^{-1} \sum_{i=1}^N \mathbf{x}_i y_i.$$

Thus, using (8) and (10), we can obtain the approximate posterior distribution of β as

$$p(\beta \mid \hat{\beta}) = \frac{\phi_p(\hat{\beta}; \beta, \hat{V}_\beta) \pi(\beta)}{\int \phi_p(\hat{\beta}; \beta, \hat{V}_\beta) \pi(\beta) d\beta}, \quad (11)$$

where ϕ_p denotes a p -dimensional multivariate normal density and $\pi(\beta)$ is a prior distribution for β .

Now, we wish to find the posterior distribution of Y for a given β . First, define

$$\hat{Y}_{\text{reg}}(\beta) = \bar{\mathbf{X}}^t \beta + \frac{1}{N} \sum_{i \in A} \frac{1}{\pi_i} (y_i - \mathbf{x}_i^t \beta),$$

Note that $\hat{Y}_{\text{reg}}(\boldsymbol{\theta})$ is a design unbiased estimator of \bar{Y} , regardless of $\boldsymbol{\theta}$. Under some regularity conditions, we can show that $\hat{Y}_{\text{reg}}(\boldsymbol{\theta})$ follows a normal distribution asymptotically. Thus, we obtain

$$\frac{\hat{Y}_{\text{reg}}(\boldsymbol{\theta}) - \bar{Y}}{\sqrt{\hat{V}_e(\boldsymbol{\theta})}} \mid \bar{Y}, \boldsymbol{\theta} \xrightarrow{\mathcal{L}} N(0, 1), \quad (12)$$

where

$$\hat{V}_e(\boldsymbol{\theta}) = \frac{1}{N^2} \sum_{i \in A} \sum_{j \in A} \frac{\Delta_{ij}}{\pi_{ij}} \frac{1}{\pi_i} \frac{1}{\pi_j} (y_i - \mathbf{x}_i^t \boldsymbol{\theta})(y_j - \mathbf{x}_j^t \boldsymbol{\theta}), \quad (13)$$

is a design consistent variance estimator of $\hat{Y}_{\text{reg}}(\boldsymbol{\theta})$ for given $\boldsymbol{\theta}$. We then use $\varphi(\hat{Y}_{\text{reg}}(\boldsymbol{\theta}); \bar{Y}, \hat{V}_e(\boldsymbol{\theta}))$ as the density for the approximate sampling distribution of $\hat{Y}_{\text{reg}}(\boldsymbol{\theta})$ in (12), where $\varphi(\cdot; \mu, \sigma^2)$ is the normal density function with mean μ and variance σ^2 . Thus, the approximate posterior distribution of \bar{Y} given $\boldsymbol{\theta}$ can be defined as

$$p(\bar{Y} \mid \hat{Y}_{\text{reg}}(\boldsymbol{\theta}), \boldsymbol{\theta}) \propto \varphi(\hat{Y}_{\text{reg}}(\boldsymbol{\theta}); \bar{Y}, \hat{V}_e(\boldsymbol{\theta})) \pi(\bar{Y} \mid \boldsymbol{\theta}), \quad (14)$$

where $\pi(\bar{Y})$ is a conditional prior distribution of \bar{Y} given $\boldsymbol{\theta}$. Without extra assumptions, we can use a flat prior distribution for $\pi(\bar{Y} \mid \boldsymbol{\theta})$. See Remark 1 below.

Therefore, combining (11) and (14), the approximate posterior distribution of \bar{Y} can be obtained as

$$p(\bar{Y} \mid \hat{Y}_{\text{reg}}(\hat{\boldsymbol{\beta}}), \hat{\boldsymbol{\beta}}) = \frac{\int \phi(\hat{Y}_{\text{reg}}(\boldsymbol{\beta}); \bar{Y}, \hat{V}_e(\boldsymbol{\beta})) \phi_p(\hat{\boldsymbol{\beta}}; \boldsymbol{\beta}, \hat{\mathbf{V}}_{\boldsymbol{\beta}}) \pi(\boldsymbol{\beta}) \pi(\bar{Y} \mid \boldsymbol{\beta}) d\boldsymbol{\beta}}{\iint \phi(\hat{Y}_{\text{reg}}(\boldsymbol{\beta}); \bar{Y}, \hat{V}_e(\boldsymbol{\beta})) \phi_p(\hat{\boldsymbol{\beta}}; \boldsymbol{\beta}, \hat{\mathbf{V}}_{\boldsymbol{\beta}}) \pi(\boldsymbol{\beta}) \pi(\bar{Y} \mid \boldsymbol{\beta}) d\boldsymbol{\beta} d\bar{Y}}. \quad (15)$$

Generating posterior samples from (15) can be easily carried out via the following two steps:

1. Generate posterior sample $\boldsymbol{\theta}^*$ of $\boldsymbol{\theta}$ from (11).

2. Generate posterior sample of Y from the conditional posterior (14) given β^* .

Based on the approximate posterior samples of Y , we can compute posterior mean as a point estimator as well as credible intervals for uncertainty quantification for Y including the variability in estimating β .

Remark 1. *If an intercept term is included in \mathbf{x}_i , that is, $\mathbf{a}^t \mathbf{x}_i = 1, \forall i \in \{1, \dots, N\}$, for some \mathbf{a} , then we have $Y = \mathbf{X}^t \beta$ and the parameter Y is completely determined from β . In this case, the posterior distribution in (15) reduces to*

$$p(\bar{Y} | \hat{Y}_{\text{reg}}(\hat{\beta}), \hat{\beta}) = \int p(\beta | \hat{\beta}) \pi(\bar{Y} | \beta) d\beta,$$

where $p(\beta | \hat{\beta})$ is defined in (11) and $\pi(\bar{Y} | \beta)$ is a degenerating distribution at $Y = \mathbf{X}^t \beta$.

The following theorem presents an asymptotic property of the proposed approximate Bayesian method.

Theorem 1. *Under the regularity conditions described in the Appendix, conditional on the full sample data,*

$$\sup_{\bar{Y} \in \Theta_Y} \left| p(\bar{Y} | \hat{Y}_{\text{reg}}(\hat{\beta}), \hat{\beta}) - \phi(\bar{Y}; \hat{Y}_{\text{reg}}, \hat{V}_e) \right| \rightarrow 0, \quad (16)$$

as $n \rightarrow \infty$ in probability, where Θ_Y is the feasible set for \bar{Y} and $p(\bar{Y} | \hat{Y}_{\text{reg}}(\hat{\beta}), \hat{\beta})$ is given in (15).

Theorem 1 is a special case of the Bernstein-von Mises theorem (van der Vaart, 2000, Section 10.2) in survey regression estimation, and its proof is given in the Appendix. According to Theorem 1, the credible interval for Y constructed from the approximated posterior distribution (15) is asymptotically equivalent to the frequentist confidence interval based on the asymptotic normality of the common

survey regression estimator. Therefore, the frequentist survey regression estimator can be formally interpreted by the Bayesian inference. The consistency of the Bayesian point estimator (e.g. posterior mean) follows directly from (16) since $V_e(\hat{\beta}) \rightarrow 0$ in probability as $n \rightarrow \infty$.

4 Approximate Bayesian method with shrinkage priors

We now consider the case when there are many auxiliary variables in applying regression estimation. When p is large, it is important to select suitable auxiliary variables that are associated with the response variable to prevent irrelevant covariates from rendering the resulting estimator inefficient. To this end, we assume that the regression model in (1) contains an intercept term. That is,

$$E(y_i | x_i) = \beta_0 + x_i^t \beta_1, \quad (17)$$

where β_0 is an intercept term.

To deal with the problem in the Bayesian way, we may define the approximate posterior distribution of \bar{Y} given both β_0 and β_1 as similar to (15). That is, we use the asymptotic distribution of the estimators $\hat{\beta}_0$ and $\hat{\beta}_1$ of β_0 and β_1 , respectively, and assign a shrinkage prior for β_1 and flat prior for β_0 . Let $\pi_\lambda(\beta_1)$ be the shrinkage prior for β_1 with a structural parameter λ which might be multivariate.

Among the several choices of shrinkage priors, we specifically consider two priors for β_1 : Laplace (Park and Casella, 2008) and horseshoe (Carvalho et al., 2009, 2010). The Laplace prior is given by $\pi_\lambda(\beta_1) \propto \exp(-\lambda \sum_{k=1}^p |\beta_k|)$, which is related to Lasso regression (Tibshirani, 1996), so that the proposed approximated Bayesian method can be seen as the Bayesian version of a survey regression estimator with Lasso (McConville et al., 2017). The horseshoe prior is a more advanced shrinkage prior with the form:

$$\pi_\lambda(\boldsymbol{\beta}_1) = \prod_{k=1}^p \int_0^\infty \phi(\beta_k; 0, \lambda^2 u_k^2) \frac{2}{\pi(1+u_k^2)} du_k, \quad (18)$$

where $\varphi(\cdot; a, b)$ denotes the normal density function with mean a and variance b . It is known that the horseshoe prior enjoys more severe shrinkage for the zero elements of $\boldsymbol{\beta}_1$ than the Laplace prior, thus allowing strong signals to remain large (Carvalho et al., 2009).

Let \hat{V}_β be the asymptotic variance-covariance matrix of $(\hat{\beta}_0, \hat{\boldsymbol{\beta}}_1^t)$. Then, under the flat prior for β_0 , the approximate posterior distribution of $\boldsymbol{\beta} = (\beta_0, \boldsymbol{\beta}_1^t)$ can be defined as

$$p_\lambda(\beta_0, \boldsymbol{\beta}_1 | \hat{\beta}_0, \hat{\boldsymbol{\beta}}_1) = \frac{\phi((\hat{\beta}_0, \hat{\boldsymbol{\beta}}_1^t); (\beta_0, \boldsymbol{\beta}_1), \hat{V}_\beta) \pi_\lambda(\boldsymbol{\beta}_1)}{\iint \phi((\hat{\beta}_0, \hat{\boldsymbol{\beta}}_1^t); (\beta_0, \boldsymbol{\beta}_1), \hat{V}_\beta) \pi_\lambda(\boldsymbol{\beta}_1) d\beta_0 d\boldsymbol{\beta}_1}. \quad (19)$$

The marginal posterior of $\boldsymbol{\beta}_1$ is given by

$$p_\lambda(\boldsymbol{\beta}_1 | \hat{\boldsymbol{\beta}}_1) = \frac{\phi(\hat{\boldsymbol{\beta}}_1; \boldsymbol{\beta}_1, \hat{V}_{\beta_{11}}) \pi_\lambda(\boldsymbol{\beta}_1)}{\int \phi(\hat{\boldsymbol{\beta}}_1; \boldsymbol{\beta}_1, \hat{V}_{\beta_{11}}) \pi_\lambda(\boldsymbol{\beta}_1) d\boldsymbol{\beta}_1}, \quad (20)$$

where $\hat{V}_{\beta_{11}}$ is the asymptotic variance-covariance matrix of $\hat{\boldsymbol{\beta}}_1$, which is a submatrix of \hat{V}_β . Under both shrinkage priors, we can derive efficient algorithms for doing posterior computations of $\boldsymbol{\beta}_1$ as well as Y . The details are provided in the Appendix. On the other hand, the conditional posterior of β_0 given $\boldsymbol{\beta}_1$ is the normal distribution with mean $\hat{\beta}_0 + \hat{V}_{\beta_{01}} \hat{V}_{\beta_{11}}^{-1} (\hat{\boldsymbol{\beta}}_1 - \boldsymbol{\beta}_1)$ and variance $\hat{V}_{\beta_{00}} - \hat{V}_{\beta_{01}} \hat{V}_{\beta_{11}}^{-1} \hat{V}_{\beta_{10}}$, where

$$\hat{V}^\beta = \begin{pmatrix} \hat{V}_{\beta_{00}} & \hat{V}_{\beta_{01}} \\ \hat{V}_{\beta_{10}} & \hat{V}_{\beta_{11}} \end{pmatrix}.$$

Thus, we can generate posterior samples of β_0 and $\boldsymbol{\beta}_1$ from (19) via Markov Chain Monte Carlo in which we iteratively sample from the marginal posterior distribution of $\boldsymbol{\beta}_1$ and conditional posterior distribution of β_0 given $\boldsymbol{\beta}_1$. Once $\boldsymbol{\beta}$ are

sampled from (19), we can use (14) to obtain the posterior distribution of Y for a given $\boldsymbol{\beta}$. Therefore, the approximate posterior distribution of Y can be obtained as

$$p_\lambda(\bar{Y} | \hat{Y}_{\text{reg}}(\hat{\boldsymbol{\beta}}), \hat{\boldsymbol{\beta}}) = \frac{\int \phi(\hat{Y}_{\text{reg}}(\boldsymbol{\beta}); \bar{Y}, \hat{V}_e(\boldsymbol{\beta})) \phi((\hat{\beta}_0, \hat{\boldsymbol{\beta}}_1^t); (\beta_0, \boldsymbol{\beta}_1), \hat{\mathbf{V}}_\beta) \pi_\lambda(\boldsymbol{\beta}_1) \pi(\bar{Y} | \boldsymbol{\beta}) d\boldsymbol{\beta}_1}{\iint \phi(\hat{Y}_{\text{reg}}(\boldsymbol{\beta}); \bar{Y}, \hat{V}_e(\boldsymbol{\beta})) \phi((\hat{\beta}_0, \hat{\boldsymbol{\beta}}_1^t); (\beta_0, \boldsymbol{\beta}_1), \hat{\mathbf{V}}_\beta) \pi_\lambda(\boldsymbol{\beta}_1) \pi(\bar{Y} | \boldsymbol{\beta}) d\boldsymbol{\beta}_1 d\bar{Y}} \quad (21)$$

where $\hat{\boldsymbol{\beta}} = (\hat{\beta}_0, \hat{\boldsymbol{\beta}}_1^t)^t$. Generating posterior samples from (21) can be easily carried out via the following two steps:

1. For a given λ , generate posterior sample $\boldsymbol{\beta}^*$ of $\boldsymbol{\beta}$ from (19).
2. Generate posterior sample of Y from the conditional posterior (14) given $\boldsymbol{\beta}^*$.

Remark 2. Let $\hat{\beta}_0^{(R)}$ and $\hat{\boldsymbol{\beta}}_1^{(R)}$ be the estimator of β_0 and $\boldsymbol{\beta}_1$ defined as

$$(\hat{\beta}_0^{(R)}, \hat{\boldsymbol{\beta}}_1^{(R)}) = \underset{\beta_0, \boldsymbol{\beta}_1}{\text{argmin}} \left\{ \sum_{i \in A} \frac{1}{\pi_i} (y_i - \beta_0 - \mathbf{x}_i^t \boldsymbol{\beta}_1)^2 + P_\lambda(\boldsymbol{\beta}_1) \right\}, \quad (22)$$

where $P(\boldsymbol{\beta}_1) = -2 \log \pi_\lambda(\boldsymbol{\beta}_1)$ is the penalty (regularization) term for $\boldsymbol{\beta}_1$ induced from prior $\pi_\lambda(\boldsymbol{\beta}_1)$. For example, the Laplace prior for $\pi_\lambda(\boldsymbol{\beta}_1)$ leads to the penalty term $P(\boldsymbol{\beta}_1) = 2\lambda \sum_{k=1}^p |\beta_k|$, in which $\hat{\boldsymbol{\beta}}_1^{(R)}$ corresponds to the regularized estimator of $\boldsymbol{\beta}_1$ used in McConville et al. (2017). Since the exponential of $-\sum_{i \in A} \pi_i^{-1} (y_i - \beta_0 - \mathbf{x}_i^t \boldsymbol{\beta}_1)^2$ is close to the approximated likelihood $\phi_p((\hat{\beta}_0, \hat{\boldsymbol{\beta}}_1^t); (\beta_0, \boldsymbol{\beta}_1), \hat{\mathbf{V}}_\beta)$ used in the approximated Bayesian method when n is large, the mode of the approximated posterior of $(\beta_0, \boldsymbol{\beta}_1)$ would be close to the frequentist estimator (22) as well.

Remark 3. By the frequent approach, λ is often called the tuning parameter and can be selected via a data-dependent procedure such as cross validation as used in McConville et al. (2017). On the other hand, in the Bayesian approach, we assign a

prior distribution on the hyperparameter parameter λ and consider integration with respect to the posterior distribution of λ , which means that uncertainty of the hyperparameter estimation can be taken into account. Specifically, we assign a gamma prior for λ^2 as the Laplace prior and a half-Cauchy prior for λ as the horseshoe prior (18). They both lead to familiar forms of full conditional posterior distributions of λ or λ^2 . The details are given in the Appendix.

As in Section 3, we obtain the following asymptotic properties of the proposed approximate Bayesian method.

Theorem 2. *Under the regularity conditions described in the Appendix, conditional on the full sample data,*

$$\sup_{\bar{Y} \in \Theta_Y} \left| p_\lambda(\bar{Y} | \hat{Y}_{\text{reg}}(\hat{\beta}), \hat{\beta}_0, \hat{\beta}_1) - \phi(\bar{Y}; \hat{Y}_{\text{reg}}(\hat{\beta}^{(R)}), \hat{V}_e(\hat{\beta}_0^{(R)}, \hat{\beta}_1^{(R)})) \right| \rightarrow 0, \quad (23)$$

as $n \rightarrow \infty$ in probability, where Θ_Y is the feasible set for \bar{Y} and $p_\lambda(\bar{Y} | \hat{Y}_{\text{reg}}(\hat{\beta}), \hat{\beta}_0, \hat{\beta}_1)$ is given in (21).

Theorem 2 ensures that the proposed approximate Bayesian method is asymptotically equivalent to the frequentist version in which β_1 is estimated by the regularized method with penalty corresponding to the shrinkage prior used in the Bayesian method. Moreover, the proposed Bayesian method can be extended to cases using general non-linear regression, as demonstrated in the next section.

5 An Extension to non-linear models

The proposed Bayesian methods can be readily extended to work with non-linear regression. Some extensions of the regression estimator to nonlinear models are also considered in Wu and Sitter (2001), Breidt et al. (2005), and Montanari and Ranalli (2005).

We consider a general working model for y_i as $E(y_i | x_i) = m(x_i; \beta) = m_i$ and $\text{Var}(y_i | x_i) = \sigma^2 a(m_i)$ for some known functions $m(\cdot; \cdot)$ and $a(\cdot)$. The model-assisted regression estimator for Y with β known is then

$$\hat{Y}_{\text{reg},m}(\beta) = \frac{1}{N} \left\{ \sum_{i=1}^N m(\mathbf{x}_i; \beta) + \sum_{i \in A} \frac{1}{\pi_i} (y_i - m(\mathbf{x}_i; \beta)) \right\}$$

and its design-consistent variance estimator is obtained by

$$\hat{V}_{e,m}(\beta) = \frac{1}{N^2} \sum_{i \in A} \sum_{j \in A} \frac{\Delta_{ij}}{\pi_{ij}} \frac{1}{\pi_i} \frac{1}{\pi_j} \{y_i - m(\mathbf{x}_i; \beta)\} \{y_j - m(\mathbf{x}_j; \beta)\}$$

which gives the approximate conditional posterior distribution of Y given β . That is, similarly to (14), we can obtain

$$p(Y | Y_{\text{reg},m}(\beta), \beta) \propto \varphi(Y_{\text{reg},m}(\beta); Y, \hat{V}_{e,m}(\beta)) \pi(Y | \beta). \quad (24)$$

To generate the posterior values of β , we first find a design-consistent estimator $\hat{\beta}$ of β . Note that a consistent estimator $\hat{\beta}$ can be obtained by solving

$$U(\beta) \equiv \sum_{i \in A} \pi_i^{-1} \{y_i - m(x_i; \beta)\} \mathbf{h}(x_i; \beta) = 0,$$

where $\mathbf{h}(x_i; \beta) = (\partial m_i / \partial \beta) / a(m_i)$. For example, for binary y_i , we may use a logistic model with $m(x_i; \beta) = \exp(x_i^t \beta) / \{1 + \exp(x_i^t \beta)\}$ and $\text{Var}(y_i) = m_i(1 - m_i)$, which leads to $\mathbf{h}(x_i; \beta) = x_i$.

Under some regularity conditions, we can establish the asymptotic normality of $\hat{\beta}$. That is,

$$\hat{V}_{\beta}^{-1/2} (\hat{\beta} - \beta) | \beta \xrightarrow{\mathcal{L}} N(0, I),$$

where

$$\hat{V}_{\beta} = \left\{ \sum_{i \in A} \frac{1}{\pi_i} \hat{\mathbf{h}}_i \hat{\mathbf{m}}(\mathbf{x}_i; \hat{\beta})^t \right\}^{-1} \left(\sum_{i \in A} \sum_{j \in A} \frac{\Delta_{ij}}{\pi_{ij}} \frac{\hat{e}_i \hat{\mathbf{h}}_i}{\pi_i} \frac{\hat{e}_j \hat{\mathbf{h}}_j^t}{\pi_j} \right) \left\{ \sum_{i \in A} \frac{1}{\pi_i} \hat{\mathbf{h}}_i \hat{\mathbf{m}}(\mathbf{x}_i; \hat{\beta})^t \right\}^{-1},$$

with $\hat{e}_i = y_i - m(x_i; \hat{\beta})$, $\hat{\mathbf{h}}_i = \mathbf{h}(x_i; \hat{\beta})$, and $m(x; \beta) = \partial m(x; \beta) / \partial \beta$. Note that $m'(x; \beta) = m_i(1 - m_i)x_i$ under a logistic model.

Thus, the posterior distribution of β given $\hat{\beta}$ can be obtained by

$$p(\beta | \hat{\beta}) \propto \varphi(\hat{\beta} | \beta, V_{\beta})\pi(\beta). \quad (25)$$

We can use a shrinkage prior $\pi(\beta)$ for β in (25) if necessary. Once β^* is generated from (25), the posterior values of Y are generated from (24) for a given β^* .

This formula enables us to define the approximate posterior distribution of β of the form (11), so that the approximate Bayesian inference for Y can be carried out in the same way as in the linear regression case. Note that Theorem 1 still holds under the general setup as long as the regularity conditions given in the Appendix are satisfied.

6 Simulation

We investigate the performance of the proposed approximate Bayesian methods against standard frequentist methods using two limited simulation studies. In the first simulation, we consider a linear regression model for a continuous y variable. In the second simulation, we consider a binary y and apply the logistic regression model for the non-linear regression estimation.

6.1 Simulation study: linear regression

In the first simulation, we generate $x_i = (x_{i1}, \dots, x_{ip})^t$, $i = 1, \dots, N$, from a multivariate normal distribution with mean vector $(1, \dots, 1)^t$ and variance-covariance matrix

$2R(0.2)$, where $p^* = 50$ and the (i,j) -th element of $R(\rho)$ is $\rho^{|i-j|}$. The response variables Y_i are generated from the following linear regression model:

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_{p^*} x_{ip^*} + \varepsilon_i, \quad i = 1, \dots, N,$$

where $N = 10,000$, $\varepsilon_i \sim N(0,2)$, $\beta_1 = 1$, $\beta_4 = -0.5$, $\beta_7 = 1$, $\beta_{10} = -0.5$ and the other β_k 's are set to zero. For the dimension of the auxiliary information, we consider four scenarios for p of 20,30,40 and 50. For each p , we assume that we can access only $(x_{i1}, \dots, x_{ip})^t$ a subset of the full information $(x_{i1}, \dots, x_{ip^*})^t$. Note that for all scenarios the auxiliary variables significantly related with Y_i are included, and so only the amount of irrelevant information gets larger as p gets larger. We consider two scenarios for the sampling probability: (A) $\pi_i = 0.04$ and (B) $\text{logit}(1 - \pi_i) = 3.1 + 0.1y_i$. The sampled units are selected via Poisson sampling, which leads to an average sample size of around 400 in both scenarios.

The parameter of interest is $\bar{Y} = N^{-1} \sum_{i=1}^N Y_i$. We assume that $\bar{X}_k = N^{-1} \sum_{i=1}^N x_{ik}$ is known for all $k = 1, \dots, p$.

For the simulated dataset, we apply the proposed approximate Bayesian methods with the uniform prior $\pi(\beta_1) \propto 1$, Laplace prior and horseshoe prior (18) for β_1 , which are denoted by AB, ABL and ABH, respectively. For all the Bayesian methods, we use $\pi(Y) \propto 1$. We generate 5,000 posterior samples of Y after discarding the first 500 samples and compute the posterior mean of Y as the point estimate. As for the frequentist methods, we apply the original generalized regression estimator without variable selection (GREG) as well as the GREG method with Lasso regularization (GREG-L; McConville et al., 2017) and ridge estimation of β_1 (GREG-R; Rao and Singh, 1997). We also apply the Horwitz-Thompson (HT) estimator as a benchmark for efficiency comparison. In GREG-L,

the tuning parameter is selected via 10-fold cross validation, and we use the gamma prior $\text{Ga}(\lambda_*^2, 1)$ for λ^2 in ABL, where λ_* is the selected value for λ in GREG-L. In ABH, we assign a prior for the tuning parameter and generate posterior samples. Based on 1,000 Monte Carlo samples, we calculate the mean squared errors (MSE), the coverage probabilities (CP) and the average length of the 95% confidence (credible) intervals, which are reported in Table 2.

Based on the results, MSEs of AB and GREG are almost the same in all cases since AB is an approximate Bayesian version of GREG. Since AB can take account of the variability in estimating β , the coverage probabilities of AB are closer to the nominal level (95%) than GREG, which is an important advantage of the proposed method. The GREG shows shorter confidence intervals with large values of p , as the variance estimator is negatively biased, and the coverage rate is lower than the nominal levels. As p gets larger, direct use of the auxiliary information makes the point estimates more inefficient as shown in Table 2, and the methods with shrinkage estimation of β such as ABH, ABL and GREG-L provide better point estimates than AB and GREG, in terms of MSEs. We note that GREG-R does not obtain much gain compared with other shrinkage methods. Comparing ABH, ABL and GREG-L, GREG-L tends to produce short confidence intervals whose coverage probabilities are smaller than the nominal level when p is large, but the proposed ABH and ABL methods produce wider credible intervals than GREG and have coverage probabilities closer to the nominal level.

6.2 *Simulation study: logistic regression*

In the second simulation study, we consider the binary case for y_i and apply the non-linear regression method discussed in Section 5. The binary response variable Y_i are generated from the following logistic regression model:

$$Y_i \sim \text{Ber}\left(\frac{\delta_i}{1 - \delta_i}\right), \quad \log\left(\frac{\delta_i}{1 - \delta_i}\right) = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip^*}, \quad i = 1, \dots, N,$$

where $\beta_0 = -1$ and the other settings are the same as the linear regression case. We again apply the same six methods based on a logistic regression model to obtain point estimates and confidence/credible intervals of the population mean $Y =$

$N^{-1} \sum_{i=1}^N Y_i$. The obtained MSE, CP and AL based on 1,000 Monte Carlo samples are reported in Table 3, which shows again the superiority of the proposed Bayesian approach to the frequentist approach in terms of uncertainty quantification.

7 Example

We applied the proposed methods to the synthetic income data available from the sae package (Molina and Marhuenda, 2015) in R language. In the dataset, the equivalized annual net income is observed for a certain number of individuals in each province of Spain. As auxiliary variables, we used four indicators of the four groupings of ages (16–24, 25–49, 50–64 and ≥ 65 denoted by $ag1, \dots, ag4$, respectively), the indicator of having Spanish nationality na , the indicators of education levels (primary education $ed1$ and post-secondary education $ed2$), and the indicators of two employment categories (employed $em1$ and unemployed $em2$). Moreover, we considered 13 interaction variables; $ag1*na$, $ag2*na$, $ag3*na$, $ag4*na$, $ag2*ed1$, $ag3*ed1$, $ag4*ed1$, $ag1*em1$, $ag2*em1$, $ag3*em1$, $ag4*em1$, $ed1*em1$ and $ed2*em1$. Here we focus on estimating average income in three provinces, Palencia, Segovia and Soria, where the number of sampled units are 72, 58 and 20, respectively. The number of non-sampled units were around 10^6 . In order to perform joint estimation and inference in the three provinces, we employed the following working model:

$$y_i = \sum_{h=1}^3 x_{0i}^{(h)} \beta_0^{(h)} + \mathbf{x}_i^t \boldsymbol{\beta}_1 + e_i \quad (26)$$

where $x_{0i}^{(h)} = 1$ if i belong to province h , where $h = 1$ for Palencia, $h = 2$ for Segovia, and $h = 3$ for Soria, and \mathbf{x}_i is the vector of auxiliary variables with dimension $p = 22$ (9 auxiliary variables and 13 interaction variables). Here y_i is the log-transformed net income and e_i is the error term.

Under the working model (26), the posterior distribution of \bar{Y}_h is

$$p\{\bar{Y}_h \mid \hat{Y}_{h,\text{reg}}(\beta_0^{(h)}, \boldsymbol{\beta}_1), \beta_0^{(h)}, \boldsymbol{\beta}_1\} \propto \phi(\hat{Y}_{h,\text{reg}}(\beta_0^{(h)}, \boldsymbol{\beta}_1) \mid \bar{Y}_h, \hat{V}_{e,h}(\boldsymbol{\beta})) \pi(\bar{Y}_h),$$

where

$$\hat{Y}_{h,\text{reg}} = \hat{\beta}_0^{(h)} + \bar{\mathbf{X}}_h^t \hat{\boldsymbol{\beta}}_1 + \frac{1}{N_h} \sum_{i \in A_h} \frac{1}{\pi_i} \left(y_i - \hat{\beta}_0^{(h)} - \mathbf{x}_i^t \hat{\boldsymbol{\beta}}_1 \right),$$

and

$$\hat{V}_{e,h}(\boldsymbol{\beta}) = \frac{1}{N_h^2} \sum_{i \in A_h} \sum_{j \in A_h} \frac{\Delta_{ij}}{\pi_{ij}} \frac{1}{\pi_i} \frac{1}{\pi_j} \left(y_i - \beta_0^{(h)} - \mathbf{x}_i^t \boldsymbol{\beta}_1 \right) \left(y_j - \beta_0^{(h)} - \mathbf{x}_j^t \boldsymbol{\beta}_1 \right).$$

Based on the above formulas, we performed the proposed approximate Bayesian methods for \bar{Y}_h for each h , and computed 95% credible intervals for the logtransformed average income with 5000 posterior samples after discarding the first

500 samples as burn-in period. We considered three types of priors for $\boldsymbol{\beta}_1$, flat, Laplace and horseshoe priors as considered in Section 6. We also calculated 95% confidence intervals of the log-transformed average income based on the two frequentist methods, GREG and GREG-L, using the working model (26). In applying GREG-L, the tuning parameter in the Lasso estimator was selected via 10 fold cross validation.

The 95% credible intervals of $\boldsymbol{\beta}_1$ based on the approximate posterior distributions under Laplace and horseshoe priors are shown in Figure 1, in which

the design-consistent and Lasso estimates of β_1 are also given. It is observed that the approximate posterior mean of β_1 shrinks the design-consistent estimates of β_1 toward 0 although exactly zero estimates are not produced as the frequentist Lasso estimator does. The Lasso estimate selects only one variable among 22 candidates, and the variable is also significant in terms of the credible interval in both two priors. Moreover, the two Bayesian methods detect one or two more variables to be significant judging from the credible intervals. Comparing the results from two priors, the horseshoe prior provides narrower credible intervals than the Laplace prior.

In Figure 2, we show the resulting credible and confidence intervals of the average income in the three provinces. It is observed that the proposed Bayesian methods, AB and ABL, tend to produce wider credible intervals than the confidence intervals of the corresponding frequencies methods, GREG and GREG-L, respectively, which is consistent to the simulation results in Section 6. We can also confirm that the credible intervals of ABH are slightly narrower than those of ABL, which would reflect the differences of interval lengths of β_1 as shown in Figure 1.

8 Concluding Remarks

We have proposed an approximate Bayesian method for survey regression estimation using a parametric regression model as the working model. The proposed Bayesian method captures the uncertainty in estimating regression parameters even when the number of the auxiliary variables is large. A main advantage of the proposed method is that it uses a shrinkage prior for regularized regression estimation, which not only provides an efficient point estimator, but also fully captures the uncertainty associated with model selection and parameter estimation via Bayesian inference. Although we only consider two popular prior distributions here, Laplace prior and the horseshoe prior, other priors, such as the

spike-and-slab prior (Ishwaran and Rao, 2005), can be considered. Further investigation regarding the choice of the shrinkage prior distributions will be an important research topic in the future.

Although our working model is parametric, the proposed approximate Bayesian method can be applied to other semiparametric models such as local polynomial model (Breidt and Opsomer, 2000), P-spline regression model (Breidt et al., 2005), or a neural network model (Montanari and Ranalli, 2005). By finding suitable prior distributions for the semiparametric models, the model complexity parameters will be determined automatically and the uncertainty will be captured in the approximate Bayesian framework. Such extensions are beyond the scope of this paper and will be topics for future research.

Acknowledgement

The first author was supported by Japan Society for the Promotion of Science KAKENHI grant number JP18K12757. The second author was supported by US National Science Foundation (MMS-1733572).

Appendix

A. Proof of (5)

From (4), we have

$$\begin{aligned} E(R_n) &= -E \left\{ (\hat{\mathbf{X}}_{\text{HT}} - \bar{\mathbf{X}}_N)^t (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_*) \right\} = -\text{tr} \left\{ \right. \\ &= - \sum_{j=1}^p V \left(\hat{x}_{\text{HT},j}, \hat{\beta}_j \right) = O(p/n). \quad \left. \text{Cov} \left(\hat{\mathbf{X}}_{\text{HT}}, \hat{\boldsymbol{\beta}} \right) \right\} \end{aligned}$$

Also, we can show that $V(R_n) = O(p/n^2)$. Therefore, using Chebychev inequality, we have $R_n = O_p(p/n)$ and result (5) follows.

B. Posterior computation

We provide the algorithm for generating the approximate posterior distribution of β_1 given in (20) with two shrinkage priors, Laplace and horseshoe (18) priors. Using the mixture representation of both priors, we get the following Gibbs sampling algorithm.

Laplace prior

We consider the mixture representation of Laplace distribution: $\beta_k|\tau_k \sim N(0, \tau_k^2)$ and $\tau_k^2 \sim \text{Exp}(\lambda^2/2)$, independently, for $k = 1, \dots, p$. For λ^2 , we consider the conjugate prior $\text{Ga}(a, b)$, where $\text{Ga}(a, b)$ is a gamma distribution with shape parameter a and rate parameter b . The full conditional distribution of β is multivariate normal with mean $A^{-1} \hat{V}_\beta^{-1} \hat{\beta}$ and variance-covariance matrix A^{-1} where $A = \hat{V}_\beta^{-1} + \mathbf{D}^{-1}$ with $\mathbf{D} = \text{diag}(\tau_1^2, \dots, \tau_p^2)$. The full conditional distribution of λ^2

is $\text{Ga}(a+p, b + \sum_{k=1}^p \tau_k^2/2)$, and $\tau_1^2, \dots, \tau_p^2$ are conditionally independent, with $1/\tau_j^2$ q^2 in the parametrization conditionally inverse-Gaussian with parameters $\mu = \lambda/\beta_j$ of the inverse-Gaussian density given by

$$f(x) = \sqrt{\frac{\lambda}{2\pi}} x^{-3/2} \exp\left\{-\frac{\lambda(x-\mu)^2}{2\mu^2 x}\right\}, \quad x > 0$$

Horseshoe prior

From (18), the prior for β_1 can be expressed as a hierarchy: $\beta_k|u_k \sim N(0, \lambda^2 u_k^2)$ and $u_k \sim \text{HC}(0, 1)$ independently for $k = 1, \dots, p$, where $\text{HC}(0, 1)$ is the standard half-Cauchy distribution. Using the hierarchical expression of the half-Cauchy dis-

tribution, we obtain the following Gibbs sampling steps. Let $A = \hat{V}_\beta^{-1} + \mathbf{B}^{-1}$, where $\mathbf{B} = \lambda^2 \text{diag}(u_1^2, \dots, u_p^2)$. The full conditional distribution of β is multivariate normal with mean $A^{-1} \hat{V}_\beta^{-1} \hat{\beta}$ and variance-covariance matrix A^{-1} . The full conditional distribution of u_k^2 and λ^2 are, respectively, give by

$$\text{IG}\left(1, \frac{1}{\xi_k} + \frac{\beta_k^2}{2\lambda^2}\right) \quad \text{and} \quad \text{IG}\left(\frac{p+1}{2}, \frac{1}{\gamma} + \frac{1}{2} \sum_{k=1}^p \frac{\beta_k^2}{u_k^2}\right),$$

where $\text{IG}(a,b)$ denotes an inverse-Gamma distribution with shape parameter a and rate parameter b . Here ξ_k and γ are additional latent variables, and their full conditional distributions are given by $\text{IG}(1, 1 + 1/\delta_k^2)$ and $\text{IG}(1, 1 + 1/\lambda^2)$, respectively.

C. A sketched proof of Theorem 1

To discuss the asymptotic properties of the approximate Bayesian method, we first assume a sequence of finite populations and samples with finite fourth moments as in Isaki and Fuller (1982). The finite population is a random sample from an unknown superpopulation model. Let \bar{Y}_* and β_* be the true values of \bar{Y} and β .

Let $B_n = (\bar{Y}_* - r_n, \bar{Y}_* + r_n)$ and C_n be a ball with centre β_* and radius $r_n \sim n^{\tau-1/2}$ for $0 < \tau < 1/2$. We make the following regularity assumptions

(C1) Assume that the sufficient conditions for the asymptotic normality of \hat{Y}_{reg} for $\bar{Y} \in B_n$ hold for the sequence of finite populations and samples.

(C2) Assume that the prior distribution $\pi(Y)$ is positive and satisfies a Lipschitz condition over its support Θ_Y ; that is, there exists $C_1 < \infty$ such that $|\pi(\theta_1) - \pi(\theta_2)| \leq C_1 |\theta_1 - \theta_2|$ for $\theta_1, \theta_2 \in \Theta_Y$.

(C3) Assume that $\hat{V}_\beta = V_\beta \{1 + o_P(1)\}$ and $(\hat{\beta} - \beta)^t \hat{V}_\beta^{-1} (\hat{\beta} - \beta) = (\hat{\beta} - \beta)^t V_\beta^{-1} (\hat{\beta} - \beta) \{1 + o_P(1)\}$ for any $\beta \in C_n$ and $n \rightarrow \infty$.

(C4) Assume that $\pi(\beta)$ is positive and finite over its support Θ_β .

Sufficient conditions for (C1) are discussed within various asymptotic structures (e.g. Binder, 1983; Pfeiffermann and Sverchkov, 2009). Conditions (C2) and (C4) are satisfied for common priors such as (multivariate) normal distribution. Condition (C3) essentially requires that the design variance estimators be consistent and meet a certain continuity condition.

Proof. Let $g(\bar{Y}, \beta) = \varphi(\tilde{Y}_{\text{reg}}(\beta); Y, \hat{V}_e(\beta))\varphi_p(\hat{\beta}; \beta, \hat{V}_\beta)\pi(\beta)$. Then, the approximated posterior distribution is given by

$$p(\bar{Y} | \hat{Y}_{\text{reg}}(\hat{\beta}), \hat{\beta}) = \frac{\int g(\bar{Y}, \beta) d\beta}{\iint g(\bar{Y}, \beta) d\beta d\bar{Y}}.$$

Note that

$$\int g(\bar{Y}, \beta) d\beta = \int_{\beta \in C_n} g(\bar{Y}, \beta) d\beta + \int_{\beta \in \mathbb{R}^p \setminus C_n} g(\bar{Y}, \beta) d\beta \quad (27) \text{ By the}$$

same argument in the proof of Theorem 1 in Wang et al. (2018), we have

$$\text{plim}_{n \rightarrow \infty} \int_{\beta \in C_n} \phi_p(\hat{\beta}; \beta, \hat{V}_\beta) d\beta = 1,$$

so the second term in (27) is $o_P(1)$. On the other hand, under condition (C3), $\varphi_p(\hat{\beta}; \beta, \hat{V}_\beta) = \varphi_p(\hat{\beta}; \beta, V_\beta)\{1 + o_P(1)\}$ as $n \rightarrow \infty$, for any $\beta \in C_n$, thereby under condition (C4),

$$\begin{aligned} \int_{\beta \in C_n} g(\bar{Y}, \beta) d\beta &= \int_{\beta \in C_n} \varphi(\tilde{Y}_{\text{reg}}(\beta); Y, \hat{V}_e(\beta))\varphi_p(\hat{\beta}; \beta, V_\beta)\pi(\beta) d\beta \\ &= \varphi(\tilde{Y}_{\text{reg}}(\beta_*); Y, \hat{V}_e(\beta_*))\pi(\beta_*)\{1 + o_P(1)\} \end{aligned}$$

as $n \rightarrow \infty$ since $V \rightarrow 0$ and $\hat{\beta} \rightarrow \beta_*$ as $n \rightarrow \infty$. Hence, we have

$$\begin{aligned}
p(\bar{Y}|\hat{Y}_{\text{reg}}(\hat{\beta}), \hat{\beta}) &= \frac{\pi(\beta_*)\phi(\hat{Y}_{\text{reg}}(\beta_*); \bar{Y}, \hat{V}_e(\beta_*))\pi(\bar{Y})\{1 + o_P(1)\}}{\pi(\beta_*) \int \phi(\hat{Y}_{\text{reg}}(\beta_*); \bar{Y}, \hat{V}_e(\beta_*))\pi(\bar{Y})d\bar{Y}\{1 + o_P(1)\}} \\
&= \frac{\pi(\bar{Y})}{\pi(\bar{Y}_*)}\phi(\hat{Y}_{\text{reg}}(\beta_*); \bar{Y}, \hat{V}_e(\beta_*))\{1 + o_P(1)\} \\
&= \phi(\hat{Y}_{\text{reg}}(\beta_*); \bar{Y}, \hat{V}_e(\beta_*))\{1 + o_P(1)\} \tag{28}
\end{aligned}$$

$$= \varphi(\hat{Y}_{\text{reg}}(\hat{\beta}); \bar{Y}, \hat{V}_e(\hat{\beta}))\{1 + o_P(1)\}, \tag{29}$$

for any $Y \in B_n$ as $n \rightarrow \infty$, where (28) follows from (C2), and (29) follows from

the properties $\hat{V}_e(\hat{\beta}) = \hat{V}_e(\beta_*)\{1 + o_P(1)\}$ and $\hat{Y}_{\text{reg}}(\hat{\beta}) = \hat{Y}_{\text{reg}}(\beta_*)\{1 + o_P(1)\}$ under (C1). Let $R_n = \{\bar{Y} \in \Theta_Y : \hat{V}_e(\hat{\beta})^{-1}(\hat{Y}_{\text{reg}}(\hat{\beta}) - \bar{Y})^2 \leq \chi_{1, \beta_0}^2\}$, where χ_{1, β_0}^2 is the upper β_0 -quantile of the chi-squared distribution with 1 degree of freedom. Then, $\text{plim}_{n \rightarrow \infty} P(R_n) = \beta_0$. Since $\hat{Y}_{\text{reg}}(\hat{\beta}) - \hat{Y}_{\text{reg}}(\beta_*) = O_p(n^{-1/2})$ and $r_n = n^{\tau-1/2}$, which is slower than $n^{-1/2}$, it holds that $\lim_{n \rightarrow \infty} P(R_n \subset B_n) = 1$. Then,

$$\lim_{n \rightarrow \infty} P\left(\int_{B_n} \phi(\hat{Y}_{\text{reg}}(\hat{\beta}); \bar{Y}, \hat{V}_e(\hat{\beta}))d\bar{Y} \geq \int_{R_n} \phi(\hat{Y}_{\text{reg}}(\hat{\beta}); \bar{Y}, \hat{V}_e(\hat{\beta}))d\bar{Y}\right) = 1,$$

which means that

$$\lim_{n \rightarrow \infty} P\left(\int_{B_n} \phi(\hat{Y}_{\text{reg}}(\hat{\beta}); \bar{Y}, \hat{V}_e(\hat{\beta}))d\bar{Y} \geq \beta_0\right) = 1$$

for any $\beta_0 \in (0, 1)$, implying

$$\text{plim}_{n \rightarrow \infty} \int_{B_n} \phi(\hat{Y}_{\text{reg}}(\hat{\beta}); \bar{Y}, \hat{V}_e(\hat{\beta}))d\bar{Y} = 1. \tag{30}$$

Then,

$$\begin{aligned}
&\sup_{\bar{Y} \in \Theta_Y} \left| p(\bar{Y}|\hat{Y}_{\text{reg}}(\hat{\beta}), \hat{\beta}) - \phi(\bar{Y}; \hat{Y}_{\text{reg}}(\hat{\beta}), \hat{V}_e(\hat{\beta})) \right| \\
&\leq \sup_{\bar{Y} \in B_n} \left| p(\bar{Y}|\hat{Y}_{\text{reg}}(\hat{\beta}), \hat{\beta}) - \phi(\bar{Y}; \hat{Y}_{\text{reg}}(\hat{\beta}), \hat{V}_e(\hat{\beta})) \right| \\
&\quad + \sup_{\bar{Y} \in \Theta_Y \setminus B_n} \left| p(\bar{Y}|\hat{Y}_{\text{reg}}(\hat{\beta}), \hat{\beta}) - \phi(\bar{Y}; \hat{Y}_{\text{reg}}(\hat{\beta}), \hat{V}_e(\hat{\beta})) \right|,
\end{aligned}$$

which are both $o_P(1)$ from (29) and (30). This completes the proof. \square

D. A sketched proof of Theorem 2

The condition (C4) given in the proof of Theorem 1 may not be satisfied for shrinkage priors. For example, the horseshoe prior (18) diverge at the origin $\beta_k = 0$. In what follows, let $\beta = (\beta_0, \beta_1^t)$ and define $\hat{\beta}$ and $\hat{\beta}^{(R)}$ in the same way. We use the

following alternative condition for the shrinkage prior $\pi_\lambda(\beta)$:

(C5) The regularized estimator $\hat{\beta}_R$ under penalty $-\log\pi_\lambda(\beta_1)$ is asymptotically

$$\sqrt{\quad} \quad (R)$$

normal, that is, $n(\hat{\beta} - \beta_*) \rightarrow N(0, \mathbf{C})$, where \mathbf{C} is a positive definite matrix and λ is appropriately chosen.

Under the Laplace prior, $\hat{\beta}^{(R)}$ is equivalent to the Lasso estimator, and the above property holds if $\lambda = o(n)$ (Knight and Fu, 2000; McConville et al., 2017). For general prior $\pi_\lambda(\beta_1)$, this condition holds if the assumption regarding the penalty term $P_\lambda(\beta_1)$ given in Fan and Li (2001) is satisfied.

Proof. It is noted that

$$\begin{aligned} & \phi_p((\hat{\beta}_0, \hat{\beta}^t); (\beta_0, \beta^t), \hat{\mathbf{V}}_\beta) \pi_\lambda(\beta_1) \\ & \propto \exp \left\{ -\frac{1}{2} (\hat{\beta} - \beta)^t \hat{\mathbf{V}}_\beta^{-1} (\hat{\beta} - \beta) + \log \pi_\lambda(\beta_1) \right\} \\ & = \exp \left\{ -\frac{1}{2} \sum_{i \in A} \frac{1}{\pi_i} (y_i - \beta_0 - \mathbf{x}_i^t \beta_1)^2 + \log \pi_\lambda(\beta_1) \right\} \{1 + o_P(1)\} \\ & = \exp \left\{ -\frac{n}{2} (\hat{\beta}^{(R)} - \beta)^t \mathbf{C}^{-1} (\hat{\beta}^{(R)} - \beta) \right\} \{1 + o_P(1)\}. \end{aligned}$$

Define

$$g(Y, \bar{\beta}) = \varphi(\tilde{Y}_{\text{reg}}(\beta); Y, \hat{V}_e(\beta)) \varphi(\hat{\beta}; \beta, \hat{V}_\beta) \pi_\lambda(\beta_1).$$

Then, it holds that

$$\int g(Y, \beta) d\beta = \varphi(\hat{Y}_{\text{reg}}(\beta_*); Y, \hat{V}_e(\beta_*))\{1 + o_P(1)\} \beta \in R_n$$

as $n \rightarrow \infty$, where R_n is a ball with center β_* and radius $O(n^{\tau-1/2})$ for $0 < \tau < 1/2$. Hence, the statement can be proved in the same way as the proof of Theorem 1 since $\varphi(\hat{Y}_{\text{reg}}(\beta_*); Y, \hat{V}_e(\beta_*)) = \varphi(\hat{Y}_{\text{reg}}(\hat{\beta}^{(R)}); Y, \hat{V}_e(\hat{\beta}^{(R)}))\{1 + o_P(1)\}$. \square

References

- Binder, D. A. (1983). On the variances of asymptotically normal estimators from complex surveys. *International Statistical Reviews* 51, 279–292.
- Breidt, F. J., G. Claeskens, and J. D. Opsomer (2005). Model-assisted estimation for complex surveys using penalised splines. *Biometrika* 92, 831–846.
- Breidt, F. J. and J. D. Opsomer (2000). Local polynomial regression estimators in survey sampling. *Annals of Statistics* 28, 403–427.
- Breidt, F. J. and J. D. Opsomer (2017). Model-assisted survey estimation with modern prediction techniques. *Statistical Science* 32, 190–205.
- Carvalho, C. M., N. G. Polson, and J. G. Scott (2009). Handling sparsity via the horseshoe. *Proceedings of the 12th International Conference on Artificial Intelligence and Statistics (AISTATS 2009)*.
- Carvalho, C. M., N. G. Polson, and J. G. Scott (2010). The horseshoe estimator for sparse signals. *Biometrika* 97, 465–480.
- Chen, J. K. T., R. L. Valliant, and M. R. Elliott (2018). Model-assisted calibration of non-probability sample survey data using adaptive LASSO. *Survey Methodology* 44, 117–144.

- Deville, J. C. and C. E. Sa rndal (1992). Calibration estimators in survey sampling. *Journal of the American Statistical Association* 87, 376–382.
- Fan, J. and R. Li (2001). Variable selection via nonconcave penalized likelihood and its Oracle properties. *Journal of the American Statistical Association* 96, 1348–1360.
- Fuller, W. A. (2002). Regression estimation for sample surveys. *Survey Methodology* 28, 5–23.
- Fuller, W. A. (2009). *Sampling Statistics*. Wiley.
- Godambe, V. P. and V. M. Joshi (1965). Admissibility and Bayes estimation in sampling finite populations, 1. *Annals of Mathematical Statistics* 36, 1707–1722.
- Hoerl, E. and R. W. Kennard (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics* 12, 55–67.
- Horvitz, D. G. and D. J. Thompson (1952). A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association* 47, 663–685.
- Isaki, C. T. and W. A. Fuller (1982). Survey design under the regression superpopulation model. *Journal of the American Statistical Association* 77, 89–96.
- Ishwaran, H. and J. S. Rao (2005). Spike and slab variable selection: Frequentist and Bayesian strategies. *Annals of Statistics* 33, 730–773+.
- Kim, J. K. and M. Park (2010). Calibration estimation in survey sampling. *International Statistical Review* 78, 21–39.
- Knight, K. and W. Fu (2000). Asymptotics for Lasso-type estimators. *Journal of Official Statistics* 28, 1356–1378.

- Little, R. J. A. (2012). Calibrated Bayes, an alternative inferential paradigm for official statistics. *Journal of Official Statistics*.
- McConville, K., F. Breidt, T. Lee, and G. Moisen (2017). Model-assisted survey regression estimation with the LASSO. *Journal of Survey Statistics and Methodology* 5, 131–158.
- Molina, I. and Y. Marhuenda (2015). sae: An R package for small area estimation. *The R Journal* 7, 81–98.
- Montanari, G. E. and M. G. Ranalli (2005). Nonparametric model calibration estimation in survey sampling. *Journal of the American Statistical Association* 100, 1429–1442.
- Park, T. and G. Casella (2008). The Bayesian Lasso. *Journal of the American Statistical Association* 103, 681–686.
- Pfeffermann, D. and M. Sverchkov (1999). Parametric and semiparametric estimation of regression models fitted to survey data. *Sankhyā, Series B* 61, 166–186.
- Pfeffermann, D. and M. Sverchkov (2009). Inference under informative sampling. *Handbook of Statistics*.
- Rao, J. and A. Singh (1997). A ridge-shrinkage method for range-restricted weight calibration in survey sampling. In *Proceedings of the Section on Survey Research Methods*, pp. 57–65. American Statistical Association.
- Saˆrndal, C. E. and Lundstrˆom (2005). *Estimation in surveys with nonresponse*. Chichester: John Wiley & Sons.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B* 58, 267–288.

van der Vaart, A. W. (2000). *Asymptotic Statistics*. New York: Cambridge University Press.

Wang, Z., J. K. Kim, and S. Yang (2018). Approximate Bayesian inference under informative sampling. *Biometrika* 105, 91–102.

Wu, C. and R. R. Sitter (2001). A model-calibration approach to using complete auxiliary information from survey data. *Journal of the American Statistical Association* 96, 185–193.

Zou, H. (2006). The adaptive Lasso and its Oracle properties. *Journal of the American Statistical Association* 101, 1418–1429.

Zou, H. and T. Hastie (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society, Series B* 67, 301–320.

Table 2: Summary of the simulation results in scenarios (A) and (B) with linear regression. All values are multiplied by 100.

Scenario (A): $\pi_i = 0.04$								
p	GREG	GREG-L	GREG-R	AB	ABL		ABH HT	
301.06	1.00	1.07	1.06	1.04			1.002.32	
40	1.01	0.98	1.02	1.01	1.01	0.99	2.32	1.012.32
	1.15	1.00	1.14	1.15	1.05	1.03	2.32	
	93.9	94.2	94.2	95.2	94.7	95.3	95.4	
	92.8	94.5	93.1	94.6	94.5	95.3	95.4	
	92.0	94.3	92.7	94.5	94.8	95.1	95.4	
	90.9	93.8	91.7	93.8	94.4	94.8	95.4	
	37.3	37.8	37.3	39.1	39.4	38.7	60.5	
	36.8	37.7	36.9	39.4	39.6	38.7	60.5	
	36.3	37.6	36.5	39.8	39.6	38.8	60.5	
	35.8	37.6	36.1	40.1	39.5	38.8	60.5	
		1.10	1.10	1.05				

	20	1.14	1.12	1.16	1.14	1.21	1.12	3.35
MSE	50	1.40	1.15	1.44	1.40	1.27	1.17	3.35
CP	20	92.8	93.0	92.5	94.4	94.0	94.2	94.4
	30	91.6	92.9	91.2	93.0	93.9	93.7	94.4
	40	89.9	93.4	89.8	92.5	94.3	93.6	94.4
	50	88.1	92.8	88.3	91.9	93.4	93.9	94.4
AL	20	39.1	39.9	39.3	41.1	42.9	41.0	70.3
	30	38.4	39.8	38.6	41.3	43.2	41.1	70.3
	40	37.7	39.8	38.1	41.5	43.2	41.1	70.3
	50	37.0	39.7	37.5	41.7	43.1	41.1	70.3

Scenario (B): $\text{logit}(1 - \pi_i) = 3.1 + 0.1y_i$								
p	GREG	GREG-L	GREG-R	AB	ABL		ABH	HT
30	1.21	1.12	1.24	1.21	1.24	1.13	3.35	
40	1.30	1.13	1.34	1.30	1.25	1.15	3.35	

Table 3: Summary of the simulation results in scenarios (A) and (B) with logistic regression. MSE values are multiplied by 10,000 and CP and AL values are multiplied by 100.

Scenario (A): $\pi_i = 0.04$								
p	GREG	GREG-L	GREG-R	AB	ABL		ABH	HT
303.76	3.56	3.74	3.71	3.66		3.53	12.4	

40 3.66 3.56 3.67 3.64 3.60 3.51 12.4 3.5512.4

3.97	3.58	3.94	3.86	3.81	3.60	12.4
93.0	94.2	94.1	95.4	95.7	95.7	94.3
92.4	93.9	92.8	94.7	96.2	95.8	94.3
91.1	94.0	92.4	94.9	96.4	96.0	94.3
89.6	94.1	91.7	95.2	96.3	96.2	94.3
7.01	7.12	7.10	7.45	7.58	7.44	13.6
6.88	7.09	6.99	7.52	7.69	7.52	13.6
6.73	7.08	6.89	7.57	7.81	7.59	13.6
6.58	7.06	6.82	7.61	7.94	7.69	13.6
	3.83	3.78	3.72			

Scenario (B): $\text{logit}(1 - \pi_i) = 3.1 + 0.1y_i$

p	GREG	GREG-L	GREG-R	AB	ABL	ABH	HT
-----	------	--------	--------	----	-----	-----	----

	20	3.69	3.58	3.69	3.67	3.64	3.56	12.4
MSE								
	50	4.01	3.62	3.94	3.91	3.80	3.63	12.4
CP	20	92.3	93.1	92.9	94.3	94.3	94.3	95.0
	30	91.5	92.8	92.1	94.3	94.3	94.2	95.0
	40	90.5	92.6	90.9	94.4	95.7	94.7	95.0
	50	89.0	92.6	90.1	93.4	95.8	95.1	95.0
AL	20	6.88	7.00	6.98	7.31	7.44	7.30	13.8
	30	6.75	6.97	6.87	7.37	7.54	7.36	13.8
	40	6.62	6.95	6.77	7.43	7.65	7.45	13.8
	50	6.47	6.93	6.69	7.47	7.78	7.54	13.8
	30	3.78	3.60	3.76	3.75	3.71	3.58	12.4
	40	3.88	3.60	3.85	3.82	3.75	3.61	12.4

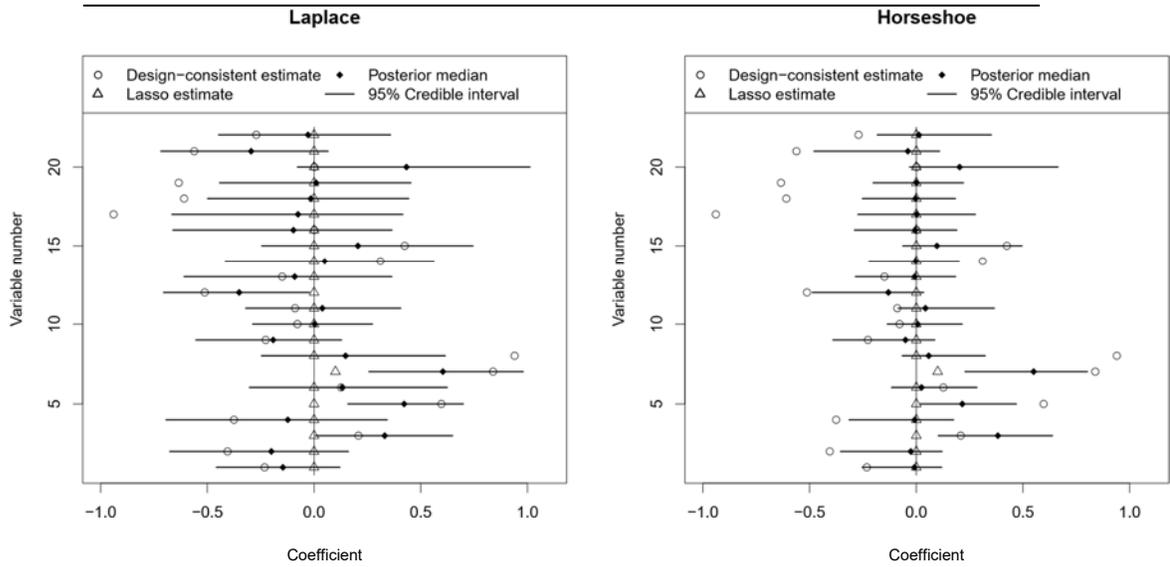


Figure 1: 95% credible intervals of regression coefficients under Laplace (left) and horseshoe (right) priors.

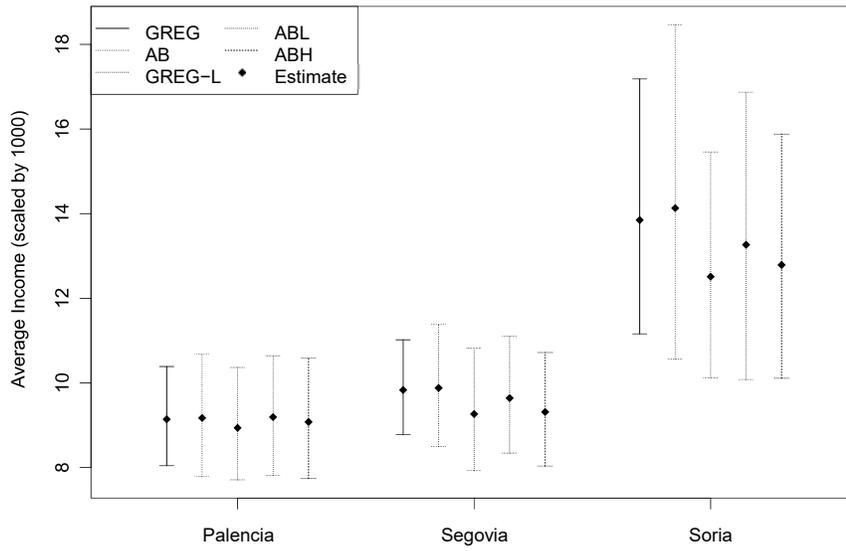


Figure 2: 95% confidence and credible intervals for average income based on five methods in three provinces in Spain.