

Spatial nonstationary hierarchical Bayes estimation of small area proportions

Priyanka Anjoy & Hukum Chandra

To cite this article: Priyanka Anjoy & Hukum Chandra (2021): Spatial nonstationary hierarchical Bayes estimation of small area proportions, Communications in Statistics - Theory and Methods, DOI: [10.1080/03610926.2021.1945632](https://doi.org/10.1080/03610926.2021.1945632)

To link to this article: <https://doi.org/10.1080/03610926.2021.1945632>



Published online: 29 Jul 2021.



Submit your article to this journal [↗](#)



Article views: 7



View related articles [↗](#)



View Crossmark data [↗](#)



Spatial nonstationary hierarchical Bayes estimation of small area proportions

Priyanka Anjoy and Hukum Chandra

ICAR-Indian Agricultural Statistics Research Institute, New Delhi, India

ABSTRACT

The hierarchical Bayes predictor of small area proportions (HBP) under an area level version of generalized linear mixed model with logit link function is widely used in small area estimation for binary variable. However, this predictor does not account for the presence of spatial nonstationarity in the data, i.e., where the parameters associated with the model covariates vary spatially. This paper develops a spatially nonstationary extension to the hierarchical Bayes predictor of small area proportions that accounts for the presence of spatial nonstationarity in the data. The proposed predictor is referred as the spatial nonstationary hierarchical Bayes predictor (HBNSP). The impact of survey design information is also explored in the proposed predictor. The empirical results from simulation studies using spatially nonstationary data indicate that the HBNSP method performs better, in terms of relative bias and relative mean squared error, than the alternative HBP method that ignore this spatial nonstationarity. The results further show that use of survey-weight to incorporate the sampling design appears to be imperative when sample data is informative. The HBNSP approach is illustrated by applying it to estimation of incidence of indebtedness in farm households across the districts in the state of Bihar in India using debt investment survey data. A map depicting the spatial distribution of incidence of indebtedness in Bihar has also been produced which provides a useful information for the government departments and ministries involved in farm credit distribution related policy planning and monitoring.

ARTICLE HISTORY

Received 17 August 2020
Accepted 15 June 2021

KEYWORDS

Spatial nonstationarity; hierarchical Bayes; binary data; indebtedness; spatial mapping

MATHEMATICAL SUBJECT CLASSIFICATION

62D05

1. Introduction

In recent years, small area estimation (SAE) technique has emerged as one of the most important topics in survey estimation because of an increasing demand for reliable small area statistics by various government and international agencies, see for example, Rao and Molina (2015). United Nations Sustainable Development Agenda has also marked the developmental strategy through availing and utilizing disaggregate level statistics in the programs and planning aimed at uprooting social and regional inequalities. Sample surveys are generally designed so that direct estimators (i.e., estimators that use only the sample data from the domain of interest) for larger domains provide reliable

CONTACT Priyanka Anjoy  anjoypriyanka90@gmail.com  ICAR-Indian Agricultural Statistics Research Institute, New Delhi 110012, India.

© 2021 Taylor & Francis Group, LLC

estimates for parameters of interest. On many occasions, however, the interest is in estimating parameters for domains that contain only a small number of sample observations or sometimes no sample observations. The term 'small areas' is used to describe domains whose sample sizes are not large enough to allow sufficiently precise direct estimation. Hereafter, we refer to these smaller domains as 'small areas' or simply 'areas'. When direct estimation is not possible, one has to rely on alternative, model-based methods for producing small area estimates. Further, large scale surveys produce reliable estimates at higher geographical level and such estimates often mask variations which is available at local levels. This restricts targeting of heterogeneity at higher levels of spatial disaggregation and also limits the scope for monitoring and evaluation of parameters locally within and across administrative units. Model-based SAE techniques are now widely used in practice to meet the indispensable need of reliable disaggregate level statistics from the existing survey data. Such SAE methods depend on the availability of population level auxiliary information related to the variable of interest, and are commonly referred to as indirect methods. The industry standard for SAE is to use unit or area level models (Fay and Herriot 1979; Battese, Harter, and Fuller 1988). In the former case these models are for the individual survey measurements and include area effects, while in the latter case these models are used to smooth out the variability in the unstable area-level direct estimates. Area-level small area modeling is usually employed when unit-level data are unavailable, or, as is often the case, where model covariates (e.g., census variables) are only available in aggregate form. In this paper we solely focus on area (or aggregated) level small area modeling.

Fay-Herriot (FH) model is one of the popular examples of aggregated level small area model. For continuous survey variable, this model is widely used in practice and has led the phenomenal development of small area literatures based on this model. However, binary or count data is often of interest in many practical applications. In epidemiological, environmental, poverty related studies such data is much common, where interest generally lies in estimation of proportions. A generalized linear mixed model (GLMM) with logit link function (also referred to as logistic linear mixed model) is commonly used for estimation of small area proportions. The basic structure of area level small area models include sampling model for direct survey estimates and associated sampling error; linking model to link the parameter of interest with area-specific auxiliary variables and random effects. The area random effect in small area models explains unstructured heterogeneity between areas which is not possible through a fixed-effect kind of structure. Two basic approaches for drawing inferences about the small area parameters of interest are known to be popular: The empirical best prediction method is based on frequentist idea to estimate unknown model parameters and the hierarchical Bayes (HB) approach assumes particular prior distributions for the hyperparameters to obtain posterior quantities of the parameter of interest. The HB approach has the flexibility to deal with complex SAE model as it overcomes the difficulties of analytical mean squared error (MSE) estimation in frequentist set up and provides quick and easier posterior variance computation based on Markov Chain Monte Carlo (MCMC) simulation. Refer Jiang and Lahiri (2001), You and Zhou (2011), Liu, Lahiri, and Kalton (2014), Rao and Molina (2015) and Chandra, Kumar, and Aditya (2018) for frequentist and Bayesian related studies and various real life applications.

This paper in particular focuses on estimation of small area proportions in hierarchical Bayes framework. Among the previous literatures, Liu, Lahiri, and Kalton (2014) and Anjoy, Chandra, and Basak (2019) have applied hierarchical Bayes version of GLMM (HBGLMM) to estimate survey-weighted small area proportions considering different cases of known and unknown sampling variance structure (denoted by HBP). The linking model of GLMM incorporates random effect which is assumed to be independent and identically distributed. As a result, spatial association-ship between geographical areas cannot be described through this structure of the model. However, in many small area problems like disease prevalence and poverty estimation, spatial contiguity between neighboring areas is very common. Therefore, induction of spatial variability in GLMM can be a way of reducing the variances or Coefficient of Variation (CV) in final estimates. One approach to incorporating such spatial dependency among the areas is to extend the GLMM to allow for spatially correlated area effects using, for example, a Simultaneous Autoregressive (SAR) model (Cressie 1993). This model allows for spatial correlation in the area effects, while keeping the fixed effects parameters spatially invariant (Chandra and Salvati 2018). There are data situations, where this assumption is inappropriate and parameters associated with the model covariates (i.e., the fixed effects parameters) vary spatially. This phenomenon is often referred to as spatial nonstationarity (Brunsdon, Fotheringham, and Charlton 2010). An alternative approach to incorporating spatial information in SAE is therefore to assume that the parameters associated with the model covariates vary spatially. In frequentist framework, Chandra, Salvati, and Chambers (2017) has devised the concept of spatial nonstationarity in area level version of GLMM (NSGLMM) for estimating small area proportions. A key feature of this approach is that it tries to capture spatial variability through incorporating spatially varying covariates in the linking model. It is worth noting that Chandra, Salvati, and Chambers (2017) approach does not use the sampling weights or clustering information in estimation of small area proportions under NSGLMM. However, use of this sampling information is essential for valid inference from survey data collected by complex survey designs. Baldermann, Salvati, and Schmid (2018) has also forwarded spatial nonstationarity concept for explaining spatial variability between areas, but their model is for unit-level data. Contrary to the previous studies, this paper describes a spatial nonstationary version of hierarchical Bayes approach for SAE that incorporates the sampling information when estimating small area proportions under an area level small area models (denoted by HBNSP). Unlike frequentist approach, the HBNSP method offers the flexibility of MSE estimation through posterior variance computation based on MCMC simulation.

Standard model-based approaches to the analysis often ignore the sampling mechanism. The GLMM technique implicitly considers equal probability sampling (simple random sampling with replacement) within each small area and thus ignores the survey-weight (Chandra, Chambers, and Salvati 2019). But, this may result in potentially large biases in the final estimates. In FH model for estimation of small area population mean, direct design-based estimators are modeled directly and the survey variance of the associated direct estimator is introduced into the model via the design-based errors. The Horvitz-Thompson estimator, weighted Hájek estimator are the structures here to incorporate survey design information (Hidiroglou and You 2016). However, this

method for continuous data requires extension for binary or count data for estimating more representative small area proportions. Consequently, the strategic idea is to modeling survey-weighted proportions (Liu, Lahiri, and Kalton 2014). Hence, the next attempt in this paper is to check the impact of complex survey design information in HBNSP. In next section we describe two version of HBNSP predictor denoted as HBNSP1 and HBNSP2 respectively taking care of HB modeling of unweighted and survey-weighted small area proportions. In Section 3, empirical evaluation studies first include a model-based simulation set up to evaluate the performance of proposed HBNSP as compared to HBP. Secondly, a design-based simulation study is carried out for comparing the performance of HBNSP1 and HBNSP2 which respectively, ignores and considers the modeling of survey-weighted proportions. As a motivating application, incidence of indebtedness is estimated across the districts in Indian state of Bihar. This application utilizes data from All India Debt-Investment Survey (AIDIS) conducted by the National Sample Survey Office (NSSO) of India for the year 2012–13 and Agriculture Census 2011–12. Estimation and spatial mapping of incidence of indebtedness pattern in farm households at disaggregate level may assist the government departments and ministries to know the liability status of farmers and thereby planning targeted policies. The paper ends with relevant concluding remarks.

2. Methodological development

Let us consider a finite population U of size N which is partitioned into D distinct small areas or simply areas. The set of population units in area i is denoted as U_i with known size N_i , such that $U = \cup_{i=1}^D U_i$ and $N = \sum_{i=1}^D N_i$. A sample s of size n is drawn from population U using a probabilistic mechanism. This resulted in sample s_i in area i with size n_i , so that $s = \cup_{i=1}^D s_i$ and $n = \sum_{i=1}^D n_i$. Assume that y_{ij} be the value of target variable y for unit j ($j=1, \dots, n_i$) in small area i . The target variable with values y_{ij} has binary response, taking value either 1 or 0. Our aim is to estimate the small domain proportions $P_i = N_i^{-1} \sum_{j=1}^{N_i} y_{ij}$. When the sample s is drawn following a complex survey design, then with each unit y_{ij} in small area i certain design weight w_{ij} is also attached, which is alternatively known as survey-weights or sampling weights.

2.1. Estimation of small area proportions

The area level version GLMM is widely used for estimation of small area proportions to improve the precision of direct survey estimates. Consider, p_i be the direct survey estimator for the parameter of interest P_i . In aggregated level model, it is customary to assume that,

$$p_i = P_i + e_i; i = 1, \dots, D$$

where e_i 's are independent sampling error associated with direct estimator p_i . Sampling error e_i is assumed to have zero mean and known sampling variance $\sigma_{e_i}^2$. The linking model of P_i attempt to relate area-specific auxiliary variables and random effect component,

$$g(P_i) = \mathbf{x}'_i \boldsymbol{\beta} + v_i; i = 1, \dots, D$$

where the linking function $g(\cdot)$ is logit for binary data and log for count data, \mathbf{x}_i represent matrix of area-specific auxiliary variables, $\boldsymbol{\beta}$ is the regression coefficient or fixed effect parameter vector and v_i being the area-specific random effect, independent and identically distributed as $E(v_i) = 0$ and $\text{var}(v_i) = \sigma_v^2$. Random area-specific effects are included in the linking model to account for between areas dissimilarities. Working under HB set up, when HBP is used for estimation of small area proportions certain prior distributions are assumed for the hyperparameters. For estimating small area proportions P_i , the sampling and linking models of HBP are represented as,

$$p_i | P_i \sim N(P_i, \sigma_{ei}^2), i = 1, \dots, D \quad \text{and} \quad \text{logit}(P_i) | \boldsymbol{\beta}, \sigma_v^2 \sim N(\mathbf{x}'_i \boldsymbol{\beta}, \sigma_v^2), i = 1, \dots, D$$

Following standard literature, prior choice for $\boldsymbol{\beta}$ is usually taken to be $N(0, \sigma_0^2)$ and for σ_v^2 prior choice is $IG(a_0, b_0)$, (IG stands for Inverse Gamma) where σ_0^2 is set to be very large (say, 10^6) and very small value for a_0 and b_0 (usually $a_0 = b_0 \rightarrow 0$) to reflect lack of prior knowledge about variance parameters (Rao and Molina 2015; You and Zhou 2011). Then, inferences about the small area parameter of interest are drawn from posterior distribution. Posterior mean is taken as the point estimate of the parameter and posterior variance as a measure of the uncertainty associated with the estimate. However, an inbuilt postulation in HBP is that fixed effect parameter or regression coefficient vector $\boldsymbol{\beta}$ is spatially invariant, this is what customarily known as spatial stationarity. In contrary, spatial nonstationarity approach tends to describe/define spatially varying regression parameters, i.e., values of the regression coefficients are necessarily different at different spatial locations. Small area estimation of proportions in presence of such spatial nonstationarity is described in next subsection.

2.2. Hierarchical Bayes version of spatial nonstationary GLMM

For spatial nonstationary version of HBGLMM or HBNSP, regression coefficients in the small area model may be expressed as explicit functions of the spatial points of the sample observations. Unlike HBP, where we restrict to a single global model with fixed parameter, HBNSP technique defines local relationships to exist between study and auxiliary variables. This approach is quite similar to the geographically weighted regression (GWR) in a multiple regression framework which takes nonstationary auxiliary variables into consideration (Brunsdon, Fotheringham, and Charlton 2010). Let, l_i be the coordinates of an arbitrarily defined spatial location (longitude and latitude) for i th small area; generally this will be its centroid. Consider, $\mathbf{l} = (l_1, \dots, l_D)'$ denoting the D component vector of such spatial locations i.e., we must have available longitude and latitude for all the D spatial locations or areas of interest. Assuming that nonstationarity is characterized by an area specific vector of fixed effects, we can write,

$$\mathbf{x}'_i \boldsymbol{\beta}(l_i) = \mathbf{x}'_i \boldsymbol{\beta} + \mathbf{x}'_i \boldsymbol{\gamma}(l_i)$$

where $\boldsymbol{\beta}(l_i) = \boldsymbol{\beta} + \boldsymbol{\gamma}(l_i)$ and $\boldsymbol{\gamma}(l_i) = (\gamma_1(l_i), \dots, \gamma_p(l_i))'$. The linking model of P_i in HBNSP attempt to relate nonstationary auxiliary variables and random effect component,

$$\text{logit}(P_i) = \mathbf{x}'_i \boldsymbol{\beta}(l_i) + v_i; i = 1, \dots, D, \quad \text{with} \quad v_i \sim N(0, \sigma_v^2)$$

Aggregating D area level models lead to the population level version of the HBNSP as

$$\mathbf{p} = \mathbf{X}\boldsymbol{\beta}(\mathbf{1}) + \mathbf{v} + \mathbf{e} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\Psi}\boldsymbol{\Theta} + \mathbf{v} + \mathbf{e}$$

where $\mathbf{p} = (p_1, \dots, p_D)'$ is the vector of direct survey estimates, $\mathbf{X} = (\mathbf{x}'_1, \dots, \mathbf{x}'_D)'$ be $D \times p$ matrix of auxiliary variates, $\boldsymbol{\beta}$ is the fixed effect parameter vector, $\mathbf{v} = (v_1, \dots, v_m)'$ is a vector of domain random effects such that $\mathbf{v} \sim N(\mathbf{0}, \sigma_v^2 \mathbf{I}_D)$, \mathbf{I}_D is the unit matrix of dimension D , $\mathbf{e} = (e_1, \dots, e_D)'$ is the vector of sampling errors with $\mathbf{e} \sim N(\mathbf{0}, \boldsymbol{\Omega})$, where $\boldsymbol{\Omega} = \text{diag}\{\sigma_{ei}^2; 1 \leq i \leq D\}$ is the matrix of design variances. $\boldsymbol{\Psi} = \{\text{diag}(\mathbf{x}'_1), \dots, \text{diag}(\mathbf{x}'_D)\}'$ is a $D \times pD$ matrix of known auxiliary data; $\boldsymbol{\Theta} = (\boldsymbol{\gamma}'(l_1), \dots, \boldsymbol{\gamma}'(l_D))'$ is a spatial Gaussian random vector of dimension $pD \times 1$ such that $E(\boldsymbol{\Theta} | \boldsymbol{\Psi}, \mathbf{1}) = \mathbf{0}$ and covariance matrix $\text{var}(\boldsymbol{\Theta} | \boldsymbol{\Psi}, \mathbf{1}) = \boldsymbol{\Sigma}_\eta = \mathbf{W} \otimes (\mathbf{c}\mathbf{c}')$, where \otimes denotes the Kronecker Product. The matrix $\mathbf{W} = 1/(1 + L(l_i, l_j))$ defines the spatial distances between sample spatial locations (l_i, l_j) , specifically distances between centroids of two locations (i, j) . In general, the only constraint on the vector \mathbf{c} is that $\boldsymbol{\Sigma}_\eta = \mathbf{W} \otimes (\mathbf{c}\mathbf{c}')$ is symmetric and non-negative definite. Following Chandra, Salvati, and Chambers (2017), we consider $\mathbf{c} = \sqrt{\eta} \mathbf{1}_p$, where $\eta \geq 0$ and $\mathbf{1}_p$ denotes the unit vector of order p . So, $\boldsymbol{\Sigma}_\eta = \eta \mathbf{W} \otimes (\mathbf{1}_p \mathbf{1}_p')$ involves non zero covariance $\text{cov}(\gamma_k(l_i), \gamma_h(l_j)) = \eta/(1 + L(l_i, l_j))$ between $\gamma_k(l_i)$ and $\gamma_h(l_j)$ for sample spatial locations (i, j) , with $k \neq h = 1, \dots, p$ and diagonal elements as η . The parameter η denotes the strength of spatial heterogeneity being explained by nonstationary auxiliary variables. In particular $\eta = 0$ indicates the situation where the model is spatially homogeneous. In HB framework, the sampling and linking models for HBNSP are then expressed as

$$\mathbf{p} | \mathbf{P} \sim N(\mathbf{P}, \boldsymbol{\Omega}) \text{ and } \text{logit}(\mathbf{P}) | \boldsymbol{\beta}, \eta, \sigma_v^2 \sim N(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Psi}\boldsymbol{\Sigma}_\eta\boldsymbol{\Psi}' + \sigma_v^2 \mathbf{I}_D)$$

The prior for hyper-parameter $\boldsymbol{\beta}$ is $N(0, \sigma_0^2)$ and for variance parameters η and σ_v^2 prior is $IG(a_0, b_0)$, where σ_0^2 is set to be very large (say, 10^6) and very small value for a_0 and b_0 . Note that HBNSP reduces to HBP when $\eta = 0$. Gibbs sampling method is implemented to estimate posterior mean $E(P_i | \mathbf{p})$ and posterior variance $\text{var}(P_i | \mathbf{p})$. The required full conditional distribution of parameters under HBP and HBNSP models are given in Section 2.3.

2.3. Survey-weighted estimation

The HB modeling of respectively unweighted and survey-weighted small area proportions is a way to check the impact of complex survey design information in the resultant estimates. It is believed that, survey-weighted direct estimates used for HB modeling purpose have the potentiality to reduce the bias or design error of the final estimates. Consider sample s of size n is drawn from population U using a complex design or at least unequal probability scheme. Let p_{ij} be the selection probability attached to j th sampling unit y_{ij} in the area i . The basic design weight can be given by $w_{ij} = (n_i p_{ij})^{-1}$. These weights can be adjusted to account for non-response and/or auxiliary information (Hidiroglou and You 2016). Normalized survey-weights d_{ij} may also be constructed, $d_{ij} = w_{ij}(\sum_j w_{ij})^{-1}$. Liu, Lahiri, and Kalton (2014) and Anjoy, Chandra, and Basak (2019) have considered HB modeling of survey-weighted small area proportions, where

GLMM structure was used for estimation of area proportions. But the effects of taking informative samples were not discussed. Here, we define two alternative models of HBNSP to study the impact of design informativeness while our aim is to estimate small area proportions in presence of spatial nonstationary auxiliary variables using the above furnished HBNSP technique. Let, $p_{i.uw}$ be the direct survey unweighted estimator for small area proportion P_i ,

$$p_{i.uw} = (n_i)^{-1} \sum_{j=1}^{n_i} y_{ij} \text{ and the variance of } p_{i.uw} \text{ is given as } \sigma_{ei.uw}^2 = n_i^{-1} P_i(1 - P_i).$$

The survey-weighted estimator denoted as, $p_{i.sw}$ and its variance is expressed as,

$$p_{i.sw} = \left(\sum_{j=1}^{n_i} w_{ij} \right)^{-1} \sum_{j=1}^{n_i} w_{ij} y_{ij} \text{ and } \sigma_{ei.sw}^2 = \left(\sum_{j=1}^{N_i} w_{ij} \right)^{-2} \left\{ \sum_{j=1}^{N_i} w_{ij} (w_{ij} - 1) (y_{ij} - P_i)^2 \right\}$$

Two HBNSP methods which are explored for the impact of complex survey design; we denote them as HBNSP1 and HBNSP2. These models are furnished below,

HBNSP1: Does not incorporate survey-weight

Sampling model: $\mathbf{p}_{uw} | \mathbf{P} \sim N(\mathbf{P}, \mathbf{\Omega}_{uw})$

Linking model: $g(\mathbf{P}) | \boldsymbol{\beta}, \eta, \sigma_v^2 \sim N(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Psi}\boldsymbol{\Sigma}_\eta \boldsymbol{\Psi}' + \sigma_v^2 \mathbf{I}_D)$

HBNSP2: Incorporate survey-weight

Sampling model: $\mathbf{p}_{sw} | \mathbf{P} \sim N(\mathbf{P}, \mathbf{\Omega}_{sw})$

Linking model: $g(\mathbf{P}) | \boldsymbol{\beta}, \eta, \sigma_v^2 \sim N(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Psi}\boldsymbol{\Sigma}_\eta \boldsymbol{\Psi}' + \sigma_v^2 \mathbf{I}_D)$

The required full conditional distributions of HBNSP1 and HBNSP2 under Gibbs sampler are given as below. Within the Gibbs sampler, particularly Metropolis-Hastings (M-H) algorithm is used for drawing random samples from full conditional distributions of posterior quantities. For HBP model, the full conditional distributions for the Gibbs sampler are given as,

$$P_i | \boldsymbol{\beta}, \sigma_v^2, p_i \propto \frac{1}{P_i(1 - P_i) \sqrt{\sigma_{ei}^2 \sigma_v^2}} \exp \left(-\frac{(p_i - P_i)^2}{2\sigma_{ei}^2} - \frac{(\log \text{it}(P_i) - \mathbf{x}'_i \boldsymbol{\beta})^2}{2\sigma_v^2} \right),$$

$$\boldsymbol{\beta} | P_i, \sigma_v^2 \sim N \left((\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \log \text{it}(\mathbf{P}), \sigma_v^2 (\mathbf{X}'\mathbf{X})^{-1} \right), \text{ and}$$

$$\sigma_v^2 | \boldsymbol{\beta}, P_i, \sim \text{IG} \left(a + \frac{D}{2}, b + \frac{\sum_{i=1}^D (\log \text{it}(P_i) - \mathbf{x}'_i \boldsymbol{\beta})^2}{2} \right)$$

For HBNSP model, the full conditional distributions for the Gibbs sampler are given as,

$$\begin{aligned} \mathbf{P}|\boldsymbol{\beta}, \eta, \sigma_v^2, \mathbf{P} &\propto |(\boldsymbol{\Psi}\boldsymbol{\Sigma}_\eta\boldsymbol{\Psi}' + \sigma_v^2\mathbf{I}_D)|^{-\frac{1}{2}}|\boldsymbol{\Omega}|^{-\frac{1}{2}}\exp\left[-\frac{1}{2}\{(\mathbf{p} - \mathbf{P})'\boldsymbol{\Omega}^{-1}(\mathbf{p} - \mathbf{P})\right. \\ &\quad \left.+ (\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta})'(\boldsymbol{\Psi}\boldsymbol{\Sigma}_\eta\boldsymbol{\Psi}' + \sigma_v^2\mathbf{I}_D)^{-1}(\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta})\}\right]\left|\frac{\partial \log \text{it}(\mathbf{P})}{\partial \mathbf{P}}\right| \\ \boldsymbol{\beta}|\mathbf{P}, \eta, \sigma_v^2 &\sim MVN\left[(\mathbf{X}'\boldsymbol{\Pi}^{-1}\mathbf{X})^{-1}(\mathbf{X}'\boldsymbol{\Pi}^{-1}\log \text{it}(\mathbf{P})), (\sigma_v^2\mathbf{I}_D + \boldsymbol{\Psi}\boldsymbol{\Sigma}_\eta\boldsymbol{\Psi}')(\mathbf{X}'\boldsymbol{\Pi}^{-1}\mathbf{X})^{-1}\right] \\ \eta|\boldsymbol{\beta}, \sigma_v^2, \mathbf{P} &\sim IG\left[a_0 + \frac{D}{2}, b_0 + \frac{(\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta} - \mathbf{v})'(\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta} - \mathbf{v})}{2}\right], \text{ and} \\ \sigma_v^2|\boldsymbol{\beta}, \eta, \mathbf{P} &\sim IG\left[a_1 + \frac{D}{2}, b_1 + \frac{(\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta} - \boldsymbol{\Psi}\boldsymbol{\Theta})'(\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta} - \boldsymbol{\Psi}\boldsymbol{\Theta})}{2}\right] \end{aligned}$$

where, $\boldsymbol{\Pi} = \boldsymbol{\Psi}\boldsymbol{\Sigma}_\eta\boldsymbol{\Psi}' + \sigma_v^2\mathbf{I}_D + \boldsymbol{\Omega}$. Recall that $\boldsymbol{\Sigma}_\eta = \eta\mathbf{W} \otimes (\mathbf{1}_p\mathbf{1}_p')$ with distance matrix $\mathbf{W} = 1/(1 + L(i, j))$

For HBNSP1 model, the full conditional distributions for the Gibbs sampler are given as,

$$\begin{aligned} \mathbf{P}|\boldsymbol{\beta}, \eta, \sigma_v^2, \mathbf{p}_{uw} &\sim |(\boldsymbol{\Psi}\boldsymbol{\Sigma}_\eta\boldsymbol{\Psi}' + \sigma_v^2\mathbf{I}_D)|^{-\frac{1}{2}}|\boldsymbol{\Omega}_{uw}|^{-\frac{1}{2}}\exp\left[-\frac{1}{2}\{(\mathbf{p}_{uw} - \mathbf{P})'\boldsymbol{\Omega}_{uw}^{-1}(\mathbf{p}_{uw} - \mathbf{P})\right. \\ &\quad \left.+ (\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta})'(\boldsymbol{\Psi}\boldsymbol{\Sigma}_\eta\boldsymbol{\Psi}' + \sigma_v^2\mathbf{I}_D)^{-1}(\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta})\}\right]\left|\frac{\partial \log \text{it}(\mathbf{P})}{\partial \mathbf{P}}\right| \\ \boldsymbol{\beta}|\mathbf{P}, \eta, \sigma_v^2 &\sim MVN\left[(\mathbf{X}'\boldsymbol{\Pi}_{uw}^{-1}\mathbf{X})^{-1}(\mathbf{X}'\boldsymbol{\Pi}_{uw}^{-1}\log \text{it}(\mathbf{P})), (\sigma_v^2\mathbf{I}_D + \boldsymbol{\Psi}\boldsymbol{\Sigma}_\eta\boldsymbol{\Psi}')(\mathbf{X}'\boldsymbol{\Pi}_{uw}^{-1}\mathbf{X})^{-1}\right] \\ \eta|\boldsymbol{\beta}, \sigma_v^2, \mathbf{P} &\sim IG\left[a_0 + \frac{D}{2}, b_0 + \frac{(\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta} - \mathbf{v})'(\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta} - \mathbf{v})}{2}\right] \text{ and} \\ \sigma_v^2|\boldsymbol{\beta}, \eta, \mathbf{P} &\sim IG\left[a_1 + \frac{D}{2}, b_1 + \frac{(\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta} - \boldsymbol{\Psi}\boldsymbol{\Theta})'(\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta} - \boldsymbol{\Psi}\boldsymbol{\Theta})}{2}\right] \end{aligned}$$

where, $\mathbf{p}_{uw} = (p_{1.uw}, \dots, p_{D.uw})'$; $\boldsymbol{\Omega}_{uw} = \text{diag}\{\sigma_{ei.uw}^2; 1 \leq i \leq D\}$ and $\boldsymbol{\Pi}_{uw} = \boldsymbol{\Psi}\boldsymbol{\Sigma}_\eta\boldsymbol{\Psi}' + \sigma_v^2\mathbf{I}_D + \boldsymbol{\Omega}_{uw}$

For HBNSP2 model, the full conditional distributions for the Gibbs sampler are given as,

$$\begin{aligned} \mathbf{P}|\boldsymbol{\beta}, \eta, \sigma_v^2, \mathbf{p}_{sw} &\sim |(\boldsymbol{\Psi}\boldsymbol{\Sigma}_\eta\boldsymbol{\Psi}' + \sigma_v^2\mathbf{I}_D)|^{-\frac{1}{2}}|\boldsymbol{\Omega}_{sw}|^{-\frac{1}{2}}\exp\left[-\frac{1}{2}\{(\mathbf{p}_{sw} - \mathbf{P})'\boldsymbol{\Omega}_{sw}^{-1}(\mathbf{p}_{sw} - \mathbf{P})\right. \\ &\quad \left.+ (\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta})'(\boldsymbol{\Psi}\boldsymbol{\Sigma}_\eta\boldsymbol{\Psi}' + \sigma_v^2\mathbf{I}_D)^{-1}(\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta})\}\right]\left|\frac{\partial \log \text{it}(\mathbf{P})}{\partial \mathbf{P}}\right| \\ \boldsymbol{\beta}|\mathbf{P}, \eta, \sigma_v^2 &\sim MVN\left[(\mathbf{X}'\boldsymbol{\Pi}_{sw}^{-1}\mathbf{X})^{-1}(\mathbf{X}'\boldsymbol{\Pi}_{sw}^{-1}\log \text{it}(\mathbf{P})), (\sigma_v^2\mathbf{I}_D + \boldsymbol{\Psi}\boldsymbol{\Sigma}_\eta\boldsymbol{\Psi}')(\mathbf{X}'\boldsymbol{\Pi}_{sw}^{-1}\mathbf{X})^{-1}\right] \\ \eta|\boldsymbol{\beta}, \sigma_v^2, \mathbf{P} &\sim IG\left[a_0 + \frac{D}{2}, b_0 + \frac{(\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta} - \mathbf{v})'(\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta} - \mathbf{v})}{2}\right], \text{ and} \\ \sigma_v^2|\boldsymbol{\beta}, \eta, \mathbf{P} &\sim IG\left[a_1 + \frac{D}{2}, b_1 + \frac{(\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta} - \boldsymbol{\Psi}\boldsymbol{\Theta})'(\log \text{it}(\mathbf{P}) - \mathbf{X}\boldsymbol{\beta} - \boldsymbol{\Psi}\boldsymbol{\Theta})}{2}\right] \end{aligned}$$

where, $\mathbf{p}_{sw} = (p_{1.sw}, \dots, p_{D.sw})'$; $\mathbf{\Omega}_{sw} = \text{diag}\{\sigma_{et.sw}^2; 1 \leq i \leq D\}$ and $\mathbf{\Pi}_{sw} = \mathbf{\Psi}\mathbf{\Sigma}_\eta\mathbf{\Psi}' + \sigma_v^2\mathbf{I}_D + \mathbf{\Omega}_{sw}$.

3. Empirical evaluations

This section reports the empirical results on the comparative performances of different estimators of the small area proportions which have been described previously. In particular, we evaluate the empirical performance of the proposed small area estimator HBNSP as compared to HBP. Further, empirical performance of HBNSP1 and HBNSP2 also has been evaluated. Two types of simulation studies are used here. Section 3.1 describes the model-based simulation set up to evaluate the performance of HBNSP and HBP. In model-based simulation, population data is generated using a specified model. In Section 3.2, a design-based simulation study is presented for comparing the performance of nonstationary process HBNSP1 and HBNSP2 which respectively, ignores and considers the modeling of survey-weighted proportions. Here, the aim is to explore impact of the incorporation of complex survey information. Simulation studies have been implemented in R. Different performance indicators considered for comparison of small area estimators are as below. Let t is the subscript for T simulations.

- $RB_i = 100 \times (T^{-1} \sum_{t=1}^T P_i^{(t)})^{-1} \left\{ T^{-1} \sum_{t=1}^T (\hat{P}_i^{(t)} - P_i^{(t)}) \right\}$ is the percentage relative bias (RB) for i^{th} small area, where $\hat{P}_i^{(t)}$ is the estimate of true population mean $P_i^{(t)}$ for i^{th} for small area at t th simulation.
- $RRMSE_i = 100 \times (T^{-1} \sum_{t=1}^T P_i^{(t)})^{-1} \left\{ \sqrt{T^{-1} \sum_{t=1}^T (\hat{P}_i^{(t)} - P_i^{(t)})^2} \right\}$ is the percentage relative root mean squared error (RRMSE) for i th for small area.
- $CR_i = 100 \times T^{-1} \sum_{t=1}^T I \left\{ LB(\hat{P}_i^{(t)}) \leq P_i^{(t)} \leq UB(\hat{P}_i^{(t)}) \right\}$ is the percentage coverage rate (CR) for i^{th} small area, where $LB(\hat{P}_i^{(t)})$ and $UB(\hat{P}_i^{(t)})$ are respectively Lower Bound (LB) and Upper Bound (UB) of the estimated population mean $\hat{P}_i^{(t)}$. $I(\cdot)$ indicates an indicator function which takes values 1 if true parameter value $P_i^{(t)}$ is within the computed interval, otherwise it takes value 0. This CR% particularly will demonstrate the credible interval property of HB models.

For design-based simulation, $P_i^{(t)}$ is equal to P_i or true population mean. A better model should show smaller values for all the performance indicators except CR. Higher the CR better is the model.

3.1. Model-based simulations

In our model-based simulations the data were generated using both stationary and non-stationary processes. In stationary data generation process (SDGP), the regression coefficients are spatially invariant. The aim of this simulation set is to examine how HBNSP performs when the data follows spatial stationary process. Here, data is generated via the linking model:

$$\text{logit}(P_i) = 1 + x_i + v_i, i = 1, \dots, D = 64$$

In case of nonstationary data generation process (NSDGP), data is generated from the following model:

$$\text{logit}(P_i) = 1 + x_i + \sqrt{\eta}(\gamma_1(l_i) + \gamma_2(l_i)x_i) + v_i, i = 1, \dots, D = 100$$

Here the values of x_i were independently drawn from the uniform distribution $x_i \sim \text{Uniform}[0, 1]$ and area random effects independently drawn as $v_i \sim N(0, \sigma_v^2 = 0.0625)$. Again the sampling model part is considered as $p_i = P_i + e_i; i = 1, \dots, D$. The independent sampling errors e_i are generated from $N(0, \sigma_{ei}^2)$ with σ_{ei}^2 taking values 0.01, 0.02, 0.03 and 0.04 respectively for equal number of areas. To define *longitude* _{i} and *latitude* _{i} of spatial locations, it is assumed that observations have been drawn from a two-dimensional grid consist of a $(\sqrt{D}x\sqrt{D})$ points uniformly spaced between -1 to 1 with a distance of $2/(\sqrt{D} - 1)$ between any two neighboring points along the vertical and horizontal axes. The D points or spatial locations are arranged in such a way that k_1 varies from -1 to 1 for each given k_2 , which also then varies from -1 to 1 . For example, when $D = 100$, the set (k_1, k_2) is, $\{k_1, k_2 = -1, -0.77, -0.55, -0.33, -0.11, 0.11, 0.33, 0.55, 0.77, 1\}$. Further, $(\gamma_1(l_i), \gamma_2(l_i))'$ has been defined as a random draw from $N(0, \mathbf{W} \otimes \mathbf{I}_2)$ with $\mathbf{W} = 1/(1 + L(l_i, l_j))$ being the distance matrix between spatial locations (l_i, l_j) . The values of η have been used as 0.5, 1, 2, 4 in this study. This simulation set up is followed from Chandra, Salvati, and Chambers (2017).

The process of generating data and estimation of small area proportions by implementing HBP and HBNSP methods was independently replicated $T = 500$ times from both stationary and nonstationary data generation process. The empirical performance and relative efficiency of the proposed HBNSP is compared with the HBP which excludes spatial nonstationarity structure. Performance of the small area HB estimators under each model is compared with respect to different prior cases for variance parameter σ_v^2 . Specifically, $IG(0.01, 0.01)$ and $IG(0.1, 0.1)$ prior cases were taken up for such sensitivity analysis with respect to prior for variance parameter σ_v^2 . However, we report the result from prior $IG(0.1, 0.1)$ only. As inferences based on different non-informative priors were found to be similar. The prior for hyperparameter β was taken as $N(0, 10^6)$. The prior for η was taken to be same as σ_v^2 . To implement the Gibbs sampler, three independent chains are used each of length 10,000. The first 5000 iterations are deleted as “burn-in” periods. Further, following Gelman and Rubin (1992), potential scale reduction factor (\hat{R}) is used to monitor the convergence of the M–H within Gibbs sampler. The \hat{R} value close to 1 is expected and equal to 1 implies stationarity.

Table 1 shows the average values of relative biases (RB), relative root mean squared errors (RRMSE) and coverage rates (CR) for HBP and HBNSP methods investigated in model-based simulations. In Table 1 these values are presented as percentage and averages over the small areas of interest ($D = 100$). Summary statistics of RB, RRMSE and CR for HBP and HBNSP methods for different values of η under NSDGP for $D = 64$ and 100 areas are reported in Appendix (Tables A1–A2). The differences between two small area predictors HBP and HBNSP in Table 1 are essentially as one would expect. When the underlying data is stationary (i.e., data generated through SDGP), with identical value of RB, value of RRMSE of HBP is marginally smaller than the HBNSP. In

Table 1. The average values of percentage relative biases (RB), percentage relative root mean squared errors (RRMSE) and percentage coverage rates (CR) for HBP and HBNSP methods in model-based simulation.

Criterion	NSDGP									
	SDGP		$\eta=0.5$		$\eta=1$		$\eta=2$		$\eta=4$	
	HBP	HBNSP	HBP	HBNSP	HBP	HBNSP	HBP	HBNSP	HBP	HBNSP
RB	0.065	0.065	-0.974	-0.478	-0.891	-0.402	-0.747	-0.359	-0.602	-0.403
RRMSE	4.289	4.447	5.636	4.635	5.130	4.242	4.593	4.047	4.622	4.281
CR	71	86	81	92	89	95	93	95	93	94

Averaged $D = 100$ areas.

contrast, in presence of nonstationarity in data (i.e., data generated through NSDGP), the HBNSP method performs consistently better than the HBP method both in terms of RB and RRMSE for all values of nonstationarity parameter η . Additionally, HBNSP has shown better coverage properties. Noncoverage rate is marginally higher for HBP method.

3.2. Design-based simulations

As in real life small area applications, we cannot be confident that our data ideally follow an assumed model, rather a working model is fitted. The endeavor of design-based simulation is to evaluate the performance of different SAE methods in the context of a realistic population. We report results from design-based simulation study based on realistic population, where a model assumption is essentially an approximation. The second type of simulation study is design-based conducted to compare the empirical performance of HBNSP1 and HBNSP2 methods. For this simulation study, debt-investment survey (AIDIS-2013) data of NSSO for rural areas of the state of Karnataka in India is used. The sample size of AIDIS-2013 is 2340 units (rural households including both indebted and non-indebted) spread over 30 districts of Karnataka. The AIDIS sample data is considered as fixed population of size 2340 units (or households) and 30 districts as small areas. Population size of small areas ranges between a minimum of 55 to a maximum of 112 with an average of 78 households. The variable of interest y_{ij} is binary which takes value 1 if a household is indebted and 0 otherwise. The aim is to estimate the proportions of indebted farm households (i.e., or incidence of indebtedness in farm households) at the district level. Here, we consider, Probability Proportional to Size with Replacement (PPSWR) samples drawn independently within each small area instead of Simple Random Sampling to take into account the effect of varying sampling weights. Motivated from the simulation set up in Hidiroglou and You (2016), PPSWR sampling was employed as follows: We first define a size measure z_{ij} for a given unit y_{ij} . Using these z_{ij} values, we computed selection probabilities $p_{ij} = z_{ij}(\sum_j z_{ij})^{-1}$ and used it to select PPSWR samples of equal size n_i from each small area. Then PPSWR samples of sizes $n_i = 10, 15, 20$ and 25 were drawn from each small area based on selection probabilities p_{ij} . This selection probability, computed from a size measure z_{ij} is a linear combination of two auxiliary variables, namely Household size and Area operated (in hectare). The basic design weight calculated as, $w_{ij} = (n_i p_{ij})^{-1}$. Further, we considered two cases for fitting HBNSP models. Case 1- No auxiliary variable is included in the

Table 2. The average values of percentage relative biases (RB), percentage relative root mean squared errors (RRMSE) and percentage coverage rates (CR) for HBNSP1 and HBNSP2 methods in design-based simulation under case 1.

Criterion	Method	$n_i = 10$	$n_i = 15$	$n_i = 20$	$n_i = 25$
RB	HBNSP2	91	96	97	97
	HBNSP1	2.52	3.28	4.15	4.10
RRMSE	HBNSP2	1.97	1.74	1.45	1.35
	HBNSP1	24.23	23.02	23.80	24.08
CR	HBNSP2	23.37	15.57	13.35	12.73
	HBNSP1	89	87	83	76

Table 3. The average values of percentage relative biases (RB), percentage relative root mean squared errors (RRMSE) and percentage coverage rates (CR) for HBNSP1 and HBNSP2 methods in design-based simulation under case 2.

Criterion	Method	$n_i = 10$	$n_i = 15$	$n_i = 20$	$n_i = 25$
RB	HBNSP2	91	95	96	97
	HBNSP1	3.45	3.77	4.90	4.88
RRMSE	HBNSP2	2.41	1.67	1.02	1.07
	HBNSP1	28.03	24.64	25.52	25.34
CR	HBNSP2	24.42	17.00	15.41	13.25
	HBNSP1	85	84	80	77

HB models and linking model contains only intercept and random effect (i.e., random mean form of model). Case 2-Available auxiliary variable (Area operated, in hectare) is used as covariate in the HB models and linking model contains intercept, one auxiliary variable and random effect (i.e., (i.e., random intercept form of model). The prior for hyperparameter β was $N(0, 10^6)$. The prior for η and σ_v^2 was taken to be $IG(0.1, 0.1)$. Gibbs sampling method is implemented with three independent chains each of length 10,000; the first 5000 iterations are deleted as “burn-in” periods. To monitor the convergence success potential scale reduction factor \hat{R} is observed. The \hat{R} value close to 1 determines that the MCMC sampler converged very well.

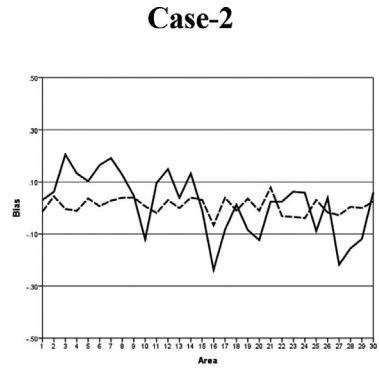
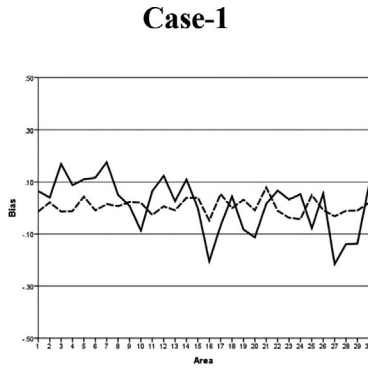
Table 2 presents the average values of RB, RRMSE and CR for the small area predictors defined by HBNSP1 and HBNSP2 methods investigated in design-based simulations under case 1. The average values of RB, RRMSE and CR for HBNSP1 and HBNSP2 under case 2 are reported in Table 3. Figure 1 plots the average values of bias for HBNSP1 and HBNSP2 methods in design-based simulations under case-1 (left side) and case-2 (right side). From these results, it is evident that design bias of survey-weighted predictor HBNSP2 is smaller than HBNSP1. Further, the values of RB for survey-weighted predictor reduces with sample size, which shows the property of design consistency of small area predictor HBNSP2. The RRMSE values are also smaller for HBNSP2 and having the same decreasing trend with increment of small area sample sizes. Investigation on coverage properties of both the models shows that noncoverage rate is higher for HBNSP1 model as compared to the other. As number of areas increases, HBNSP2 shows the better coverage percentage.

4. Application to real survey data

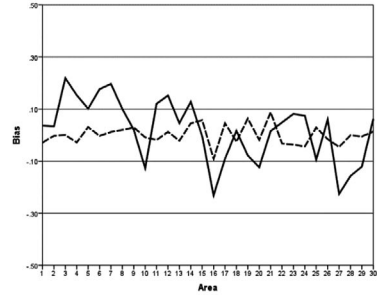
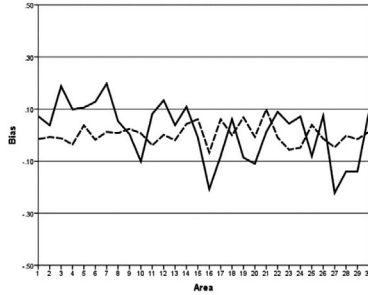
The SAE method is a cost effective and proficient approach for generating reliable micro level statistics using the existing survey data combining with auxiliary

Sample
size(n_i)

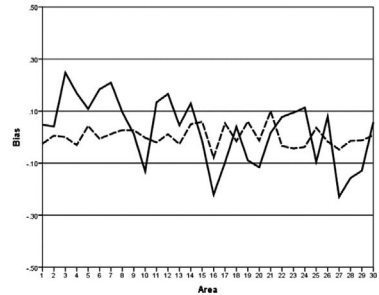
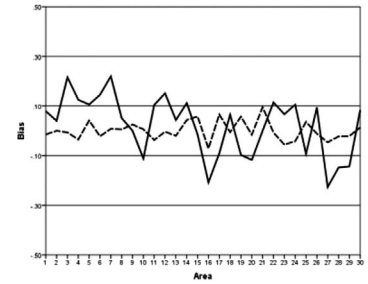
10



15



20



25

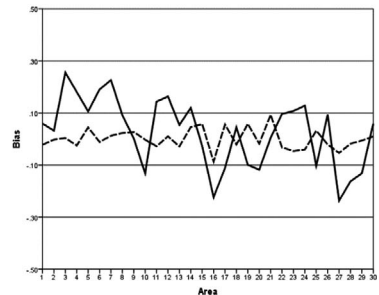
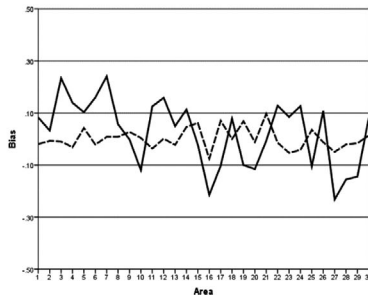


Figure 1. Comparison of bias of HBNSP1 and HBNSP2 (HBNSP1: Solid line, HBNSP2: Dash line) under case 1 (left side) and case 2 (right side).

information from other sources. This paper delineates suitable HB SAE method for generating representative and precise estimates of small area proportions in presence of nonstationarity in the data. The impact of survey design information in the proposed method for small area proportions is also explored. Empirical evaluation in the previous section has shown that, when the small area estimates are generated by incorporating spatial information in the HB models, they are more efficient than the one generated by

Table 4. Defining district categories based on AIDIS sample sizes for Bihar.

District categories	Sample sizes	District name
Small districts (S1)	52,56,84	Sheohar, Supaul, Kishanganj, Madhepura, Saharsa, Khagaria, Banka, Munger, Lakhisarai, Sheikhpura, Bhojpur, Buxar, Kaimur(Bhabua), Aurangabad, Gaya, Jamui, Jehanabad and Arwa
Medium districts (S2)	110,111,112	Pashchim Champaran, Sitamarhi, Araria, Purnia, Katihar, Gopalganj, Siwan, Vaishali, Begusarai, Bhagalpur, Nalanda, Patna and Rohtas
Large districts (S3)	126,138,139,140	Purba Champaran, Madhubani, Darbhanga, Muzaffarpur, Saran, Samastipur and Nawada

ignoring this information. Comparison of the performance of HBNSP1 and HBNSP2 through design-based simulation confirms that incorporation of survey-weight is necessary if data is obtained under informative sampling process. Here, we provide applications of the proposed models.

We now illustrate the application of HBP and HBNSP methods to estimate the proportion of indebted farm households in rural areas across the districts of Bihar using the data from AIDIS of NSSO and Agriculture Census. The AIDIS is conducted by the NSSO, Government of India at decennial intervals through household interviews from a random, nationally representative sample of households. The survey is aimed at generating average value of assets, outstanding debt per household and incidence of indebtedness, separately for the rural and urban sectors of the country, for States and Union Territories, and for different socio-economic groups. The sampling design used in this survey is stratified multi-stage random sampling with districts as strata, the census villages in the rural sector as first stage units and households as the ultimate stage units. This survey is designed and conducted to produce reliable estimates at macro or higher geographical (e.g., national and state) level. Due to small sample size limitation, this survey data cannot be used directly to generate reliable micro or local (e.g., district or further disaggregation) level estimates using traditional direct survey estimation methods. The debt investment survey data of NSSO for rural areas of Bihar state in India comprise of 3671 households (i.e., number of surveyed households from rural areas which includes both indebted and non-indebted households) spread over 38 districts. Here, different districts of the state are considered as small areas. Sample sizes in each small area vary from a minimum of 52 to maximum of 140 with an average of 97. In Table 4, districts are grouped into Small (S1), Medium (S2) and Large (S3) categories based on AIDIS sample sizes. Figure 2 presents a pictorial for the district-wise sample size distribution. The target variable y at the unit level (i.e., household) is a binary variable indicating whether a farm household is indebted or not. Quantity of interest is district-wise estimates of the incidence of indebtedness. The auxiliary variable used for HB models is Average holding size (in hectare) at district level were available from Agriculture Census data. Before we present detailed application results, it is necessary to explore whether the described data set exhibit spatial nonstationarity or not. For this purpose, district specific regression coefficients are computed by fitting GWR model. In the fitted model we have one covariate, therefore two regression coefficients (i.e., intercepts and one slope parameter).

Table 5 reports the estimated regression coefficients from GWR fit. In this table, X represents the auxiliary variable used in the application. Table 5 confirms the variation

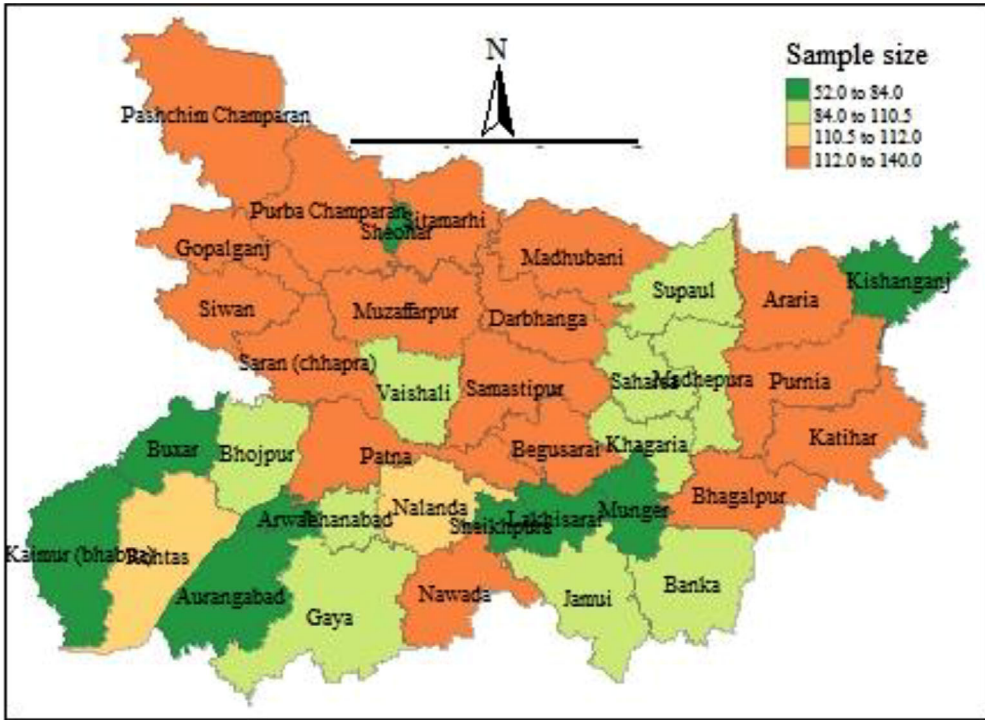


Figure 2. Map showing distribution of district-specific sample sizes in AIDIS data for Bihar.

Table 5. Summary statistics for GWR parameter estimates.

Values	Intercept	X
Minimum	-0.30	-3.56
Q1	0.18	-0.75
Mean	0.51	-0.29
Median	0.49	-0.29
Q3	0.75	0.47
Maximum	1.57	1.67

in the district specific regression coefficients which is generated through GWR model fitting. Further, Figure 3 presents the surface plot of estimated regression coefficients from GWR fit. This contour map also confirms the evidence of spatial nonstationarity in the AIDIS data. Hence, it is well expected that better performance of the small area estimates can be obtained through using spatial nonstationary HBNSP method over the non-spatial alternative HBP.

For computing HB estimates for incidence of indebtedness (P_i), we have considered prior for η and σ_v^2 as $IG(0.1, 0.1)$ and distribution of β has been taken to be $N(0, 10^6)$. The value of \hat{R} for each of the district was found to be close to 1, which implies the convergence success of MCMC sampler in each HB method implementation. Table 6 presents the summary of percentage coefficient of variation (CV%) for HBP and HBNSP estimates over different district categories. Improvement in average precision level in HBNSP method is noticeable from this Table. For small districts (S1) this improvement is much significant. Overall, this application demonstrates that HBNSP method can suitably be preferred for providing precise estimates of small area

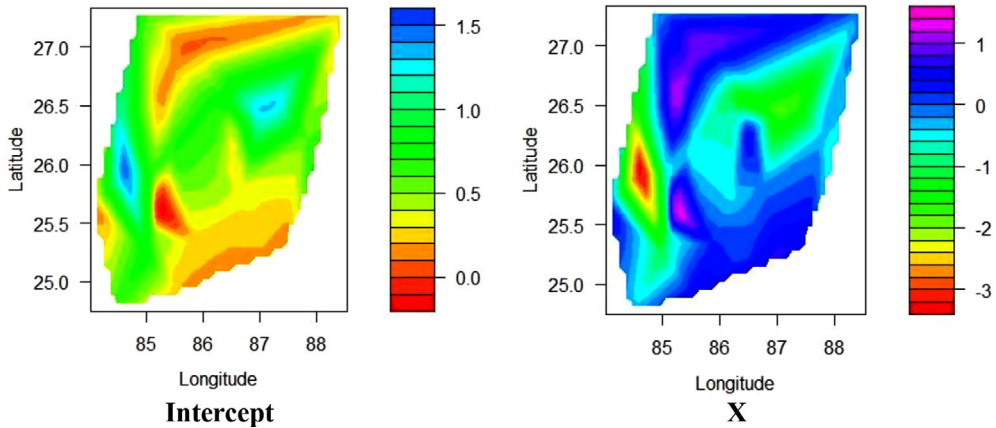


Figure 3. Contour maps showing the spatial variation in the district specific regression coefficients generated through GWR model fitting to the AIDIS data for Bihar.

proportions in presence of spatial nonstationarity in the data. The result presenting the district-wise estimates of incidence of indebtedness along with 95% confidence (credible) interval generated by HBNSP method is furnished in Appendix (Table A3).

We further illustrate the same application in AIDIS data set of Bihar in the light of showing the possible impact of incorporating survey-weight in spatial nonstationary HBNSP structures. As already discussed, the parameter of interest is the incidence of indebtedness. The survey-weights w_{ij} associated with the variable of interest y_{ij} were available from survey data file have been used. Table 7 presents the comparison of bias and relative error (RE) of HBNSP1 and HBNSP2 estimates of district-wise incidence of indebtedness. We have computed bias as the difference between average value of direct and model estimates and RE as average relative difference between direct and model estimates, see Chandra, Chambers, and Salvati (2019). Thus,

$$Bias = D^{-1} \left(\sum_i Direct\ estimate_i \right) - D^{-1} \left(\sum_i Model\ based\ estimate_i \right) \text{ and}$$

$$RE = D^{-1} \sum_i [(Direct\ estimate_i - Model\ based\ estimate_i) / Direct\ estimate_i]$$

The results in Table 7 clearly indicate that survey-weighted predictor HBNSP2 is able to provide more stable estimates than HBNSP1.

Finally, Figure 4 presents the spatial map showing distribution of proportion or incidence of indebtedness across the districts of Bihar state in India generated by the HBNSP method. Figure 5 is the spatial map of district-wise CV generated by HBNSP approach. Spatial map produced from the model-based HBNSP estimates of proportions of indebted households presents a quick view to the regional variations or disparity in district level debt/indebtedness status. Such spatial maps are certainly useful to the policy makers to frame targeted policy plans eyeing to the upliftment of deprived regions of the population or needy cultivators who are particularly indebted. In recent years, Government of India has launched number of schemes for the benefit of farmers in the country and also various loan waiver programs being run by the state administration.

Table 6. Summary statistics of percentage coefficient of variation generated by HBP and HBNSP methods.

District categories	Method	Minimum	Q1	Mean	Median	Q3	Maximum
Overall	HBNSP	7.80	8.58	9.88	9.06	11.17	12.84
	HBP	9.07	12.14	21.27	14.12	18.14	72.33
S1	HBNSP	7.14	9.10	11.01	11.15	12.74	16.80
	HBP	9.97	13.32	25.26	16.23	19.74	72.33
S2	HBNSP	7.60	11.41	12.40	12.59	13.62	16.80
	HBP	9.07	11.02	20.14	12.51	17.43	67.25
S3	HBNSP	7.14	8.58	9.71	9.51	11.09	12.58
	HBP	9.50	12.20	13.07	12.70	14.16	16.56

Table 7. Comparison of Bias and RE of HBNSP1 and HBNSP2 estimates for district-wise proportions of indebted households in Bihar.

District Categories	Bias		RE	
	HBNSP1	HBNSP2	HBNSP1	HBNSP2
Overall	-0.129	0.014	-0.934	0.024
S1	-0.054	0.017	-0.541	0.041
S2	-0.119	0.017	-0.364	0.071
S3	-0.344	0.001	-1.552	0.016

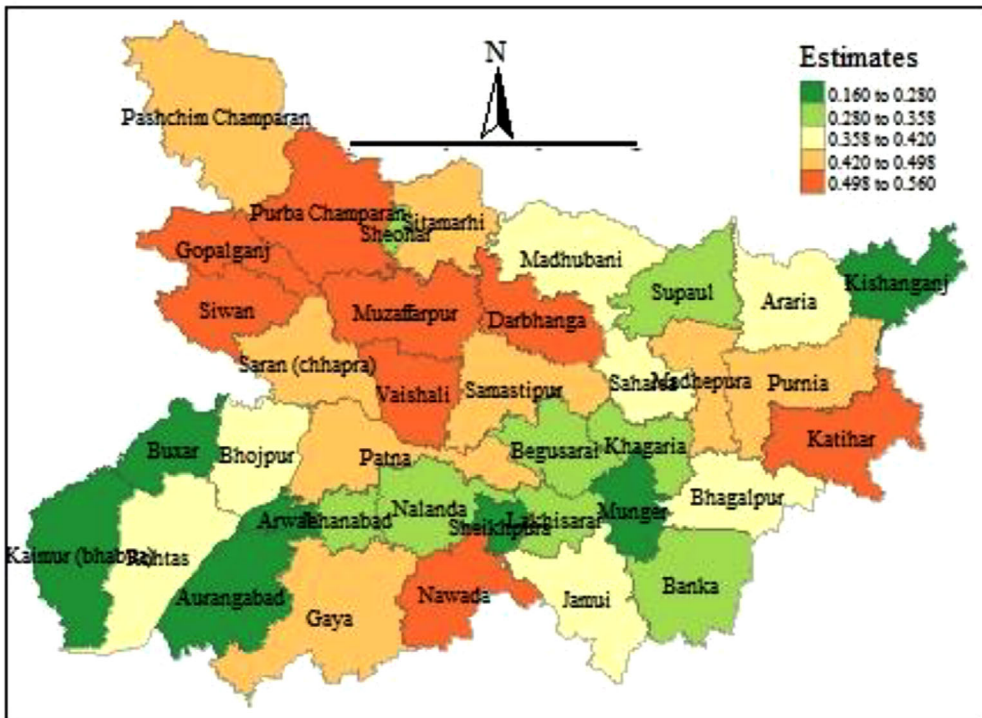


Figure 4. Spatial map of district-wise incidence of indebtedness in Bihar.

These maps will be useful to them (the administration) in identifying certain regions requiring greater level of attention.

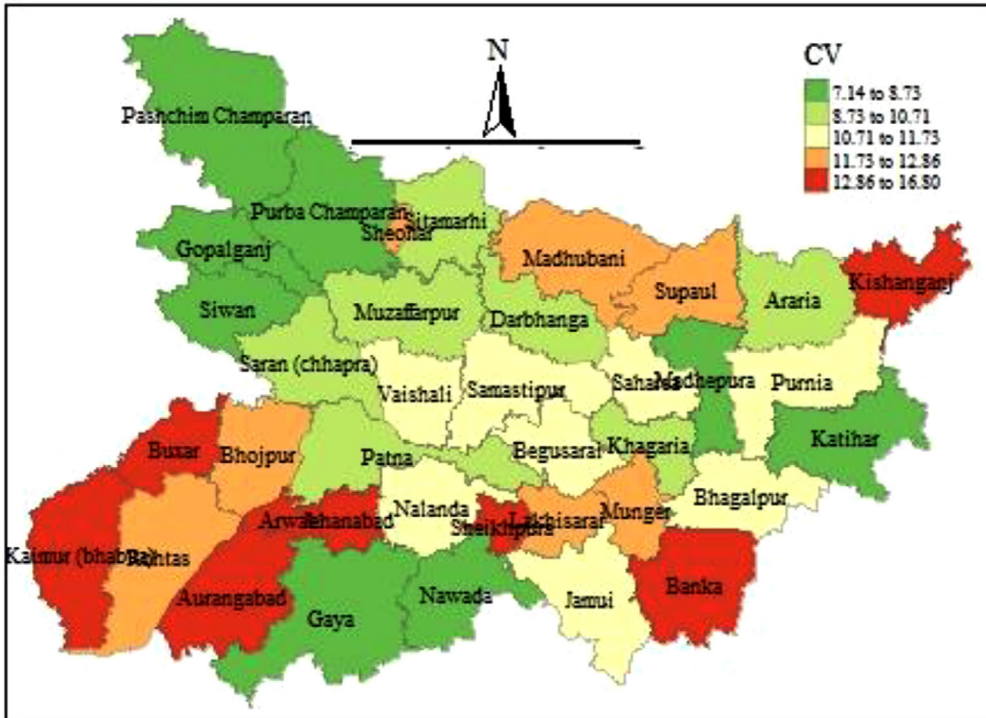


Figure 5. Spatial map showing district-wise percentage coefficient of variation for indebtedness estimates in Bihar.

5. Concluding remarks

This paper describes a spatial nonstationary extension of the area level version of the hierarchical Bayes generalized linear mixed model and considers SAE of proportions under this model. The corresponding predictor is referred to as the spatial nonstationary hierarchical Bayes predictor (HBNSP) for small area proportions. This predictor can account for the presence of spatial nonstationarity in the data where the parameters associated with the model covariates vary spatially.

Empirical results based on simulation studies provide evidence that the proposed HBNSP predictor is more efficient than the alternative hierarchical Bayes predictor under the area level generalized linear mixed model when there is a spatial nonstationarity in the data. The MSE estimation of the HBNSP predictor derived from associated posterior variance also performed reasonably well, with good coverage performance for nominal confidence intervals based on it. It is worth noting that in this paper empirical studies were also carried out using survey weights to incorporate the sampling design in SAE. This seems more realistic to implement survey weighted estimation than assuming that the sampling design is customary non-informative. The HBNSP method is also applied to real debt investment survey data to estimate the incidence of indebtedness in farm households for the districts of rural areas of the State of Bihar in India and produced a map of these districts based on these estimates of incidence of indebtedness. These estimates of incidence of indebtedness and their spatial distribution will be useful for various Government Departments and Ministries in India as well as International

organizations for their policy research and strategic planning, budget allocation and intervention for credit distribution to agricultural households.

The Census in India, like in other countries, usually has limited scope in collection of data. It focuses mainly on basic social and demographic information and that too at decennial interval. On the other hand, NSSO conducts regular surveys on a number of socio-economic indicators, but their utility is restricted to generate national and state level estimates, not administrative units below state because of small sample sizes for such units. Due to emphasis on disaggregate level Sustainable Development Goal indicators, Government of India as well as different State Governments are now struggling with generation of disaggregated level statistics. The SAE is only indispensable alternative to meet the growing demand for such disaggregated level statistics needed for decentralized policy planning. The developed SAE methodology and illustrated in the application presented in this paper can be used for calculating disaggregate level estimates of prevalence and proportions which is common in most of the socio-economic and health surveys.

Acknowledgments

The authors would like to acknowledge the valuable comments and suggestions of the Editor and an anonymous referee. These led to a considerable improvement in the paper. I also would like to deeply acknowledge the efforts made by Late Dr. Hukum Chandra in finalizing the article. He has left for heavenly abode before the final acceptance of this article.

References

- Anjoy, P., H. Chandra, and P. Basak. 2019. Estimation of disaggregate-level poverty incidence in Odisha under area-level hierarchical Bayes small area model. *Social Indicators Research* 144 (1):251–73. doi:10.1007/s11205-018-2050-9.
- Baldermann, C., N. Salvati, and T. Schmid. 2018. Robust small area estimation under spatial non-stationarity. *International Statistical Review* 86 (1):136–59. doi:10.1111/insr.12245.
- Battese, G. E., R. M. Harter, and W. A. Fuller. 1988. An error-component model for prediction of country crop areas using survey and satellite data. *Journal of the American Statistical Association* 83 (401):28–36. doi:10.1080/01621459.1988.10478561.
- Brunsdon, C., A. S. Fotheringham, and M. E. Charlton. 2010. Geographically weighted regression: a method for exploring spatial nonstationarity. *Geographical Analysis* 28 (4):281–98. doi:10.1111/j.1538-4632.1996.tb00936.x.
- Chandra, H., R. Chambers, and N. Salvati. 2019. Small area estimation of survey weighted counts under aggregated level spatial model. *Survey Methodology* 45:31–59.
- Chandra, H., S. Kumar, and K. Aditya. 2018. Small area estimation of proportions with different levels of auxiliary data. *Biometrical Journal. Biometrische Zeitschrift* 60 (2):395–415. doi:10.1002/bimj.201600128.
- Chandra, H., and N. Salvati. 2018. Small area estimation for count data under a spatial dependent aggregated level random effects model. *Communications in Statistics - Theory and Methods* 47 (5):1234–55. doi:10.1080/03610926.2017.1317806.
- Chandra, H., N. Salvati, and R. Chambers. 2017. Small area prediction of counts under a non-stationary spatial model. *Spatial Statistics* 20:30–56. doi:10.1016/j.spasta.2017.01.004.
- Cressie, N. 1993. *Statistics for spatial data*. New York: Wiley.
- Fay, R. E., and R. A. Herriot. 1979. Estimates of income for small places: An application of James-Stein procedures to census data. *Journal of the American Statistical Association* 74 (366a):269–77. doi:10.1080/01621459.1979.10482505.

- Gelman, A., and D. Rubin. 1992. Inference from iterative simulation using multiple sequences. *Statistical Science* 7 (4):457–511. doi:10.1214/ss/1177011136.
- Hidiroglou, M. A., and Y. You. 2016. Comparison of unit level and area level small area estimators. *Survey Methodology* 42:41–61.
- Jiang, J., and P. Lahiri. 2001. Empirical best prediction for small area inference with binary data. *Annals of the Institute of Statistical Mathematics* 53 (2):217–43. doi:10.1023/A:1012410420337.
- Liu, B., P. Lahiri, and G. Kalton. 2014. Hierarchical Bayes modeling of survey-weighted small area proportions. *Survey Methodology* 40:1–13.
- Rao, J. N. K., and I. Molina. 2015. *Small Area Estimation*. 2nd ed. New York: John Wiley & Sons
- You, Y., and M. Q. Zhou. 2011. Hierarchical Bayes small area estimation under a spatial model with application to health survey data. *Survey Methodology* 37:25–37.

Appendix

Table A1. Summary statistics of percentage relative biases (RB), percentage relative root mean squared errors (RRMSE) and percentage coverage rates (CR) for HBP and HBNSP methods in model-based simulations for different values of η under spatial nonstationary data generation process for $D = 100$ small areas.

Criterion	RB		RRMSE		CR	
	HBP	HBNSP	HBP	HBNSP	HBP	HBNSP
			$\eta=0.5$			
Minimum	−9.50	−8.26	2.11	1.88	22	39
Q1	−3.70	−2.46	4.28	3.41	76	91
Mean	−0.97	−0.48	5.64	4.63	81	92
Median	−0.83	−0.50	5.39	4.38	86	95
Q3	0.76	1.09	6.46	5.47	93	98
Maximum	22.30	17.34	23.72	18.52	100	100
			$\eta=1$			
Minimum	−9.30	−8.01	1.60	1.477	33	41
Q1	−3.48	−2.17	3.53	2.86	85	95
Mean	−0.89	−0.40	5.13	4.24	89	95
Median	−0.86	−0.32	4.72	4.03	94	98
Q3	0.79	1.00	6.06	5.16	98	99
Maximum	23.08	17.11	25.24	18.30	100	100
			$\eta=2$			
Minimum	−8.73	−7.63	1.36	1.22	39	44
Q1	−3.25	−2.21	2.93	2.62	91	95
Mean	−0.75	−0.36	4.59	4.05	93	95
Median	−0.58	−0.06	4.25	3.73	98	99
Q3	0.79	0.89	5.63	5.05	99	100
Maximum	23.17	19.33	25.13	20.66	100	100
			$\eta=4$			
Minimum	−8.17	−7.57	1.28	1.15	43	49
Q1	−3.14	−2.56	2.90	2.69	91	94
Mean	−0.60	−0.40	4.62	4.28	93	94
Median	−0.35	−0.24	4.27	3.75	98	99
Q3	0.78	0.83	5.72	5.28	99	99
Maximum	37.61	34.28	39.94	36.31	100	100

Table A2. Summary statistics of percentage relative biases (RB), percentage relative root mean squared errors (RRMSE) and percentage coverage rates (CR) for HBP and HBNSP methods in model-based simulations for different values of η under spatial nonstationary data generation process for $D = 64$ small areas.

Criterion	RB		RRMSE		CR	
	HBP	HBNSP	HBP	HBNSP	HBP	HBNSP
$\eta=0.5$						
Minimum	-8.21	-6.54	1.61	1.62	64	79
Q1	-2.16	-1.52	3.19	3.12	97	98
Mean	-0.34	0.00	4.18	4.10	97	98
Median	0.23	0.08	4.32	4.05	99	99
Q3	1.34	1.86	4.86	4.71	100	100
Maximum	5.76	5.66	8.63	7.13	100	100
$\eta=1$						
Minimum	-7.79	-6.35	1.39	1.49	68	81
Q1	-1.99	-1.45	2.98	2.37	98	99
Mean	-0.26	0.08	3.91	3.83	98	98
Median	0.35	0.29	4.05	3.83	100	99
Q3	1.33	1.86	4.61	4.61	100	100
Maximum	5.30	5.19	8.19	7.29	100	100
$\eta=2$						
Minimum	-7.35	-6.25	1.42	1.45	67	77
Q1	-1.76	-1.29	2.87	2.09	99	99
Mean	-0.20	0.12	3.86	3.83	98	98
Median	0.37	0.28	4.08	3.69	100	100
Q3	1.26	1.73	4.72	4.54	100	100
Maximum	4.92	4.54	7.77	7.28	100	100
$\eta=4$						
Minimum	-7.02	-5.88	1.58	1.55	68	72
Q1	-1.58	-1.09	3.26	3.42	97	98
Mean	0.30	0.28	4.40	4.36	98	98
Median	0.40	0.36	4.32	4.11	99	99
Q3	1.39	1.44	5.24	4.94	100	100
Maximum	6.25	5.80	9.48	9.45	100	100

Table A3. District-wise estimates of incidence of indebtedness along with 95% credible interval for HBNSP method of SAE for Bihar.

Districts	Sample sizes	Estimates	Lower	Upper	Districts	Sample sizes	Estimates	Lower	Upper
Araria	112	0.38	0.31	0.45	Madhubani	140	0.37	0.28	0.47
Arwal	56	0.26	0.19	0.34	Munger	56	0.27	0.21	0.33
Aurangabad	56	0.22	0.17	0.28	Muzaffarpur	138	0.51	0.42	0.6
Banka	84	0.3	0.22	0.38	Nalanda	111	0.3	0.23	0.36
Begusarai	112	0.35	0.27	0.42	Nawada	140	0.56	0.47	0.65
Bhagalpur	112	0.41	0.32	0.49	Pashchim Champaran	112	0.48	0.41	0.54
Bhojpur	84	0.37	0.28	0.46	Patna	112	0.45	0.36	0.53
Buxar	52	0.16	0.11	0.21	Purba Champaran	126	0.53	0.45	0.61
Darbhangha	140	0.56	0.46	0.66	Purnia	112	0.45	0.35	0.55
Gaya	84	0.42	0.35	0.48	Rohtas	111	0.36	0.27	0.45
Gopalganj	112	0.53	0.45	0.6	Saharsa	84	0.37	0.29	0.45
Jamui	84	0.38	0.3	0.47	Samastipur	140	0.47	0.36	0.58
Jehanabad	84	0.28	0.21	0.35	Saran	139	0.45	0.35	0.54
Kaimur	56	0.24	0.17	0.3	Sheikhpura	56	0.23	0.17	0.29
Katihar	112	0.54	0.46	0.62	Sheohar	56	0.28	0.21	0.35
Khagaria	84	0.31	0.25	0.37	Sitamarhi	112	0.43	0.35	0.51
Kishanganj	56	0.21	0.15	0.27	Siwan	112	0.56	0.47	0.65
Lakhisarai	56	0.29	0.21	0.36	Supaul	84	0.33	0.25	0.41
Madhepura	84	0.42	0.35	0.48	Vaishali	110	0.51	0.4	0.61