

1 Small area estimation under a semi-parametric covariate measured with 2 error

3 Reyhane Sefidkar^{1,2}, Mahmoud Torabi^{3*} and Amir Kavousi²

4 *Shahid Sadoughi University of Medical Sciences, Shahid Beheshti University of*
5 *Medical Sciences and University of Manitoba*

Summary

In recent years, small area estimation has played an important role in statistics as it deals with the problem of obtaining reliable estimates for parameters of interest in areas with small or even zero sample sizes corresponding to population sizes. Nested error linear regression models are often used in small area estimation assuming that the covariates are measured without error and also the relationship between covariates and response variable is linear. Small area models have also been extended to the case in which a linear relationship may not hold, using penalised spline (P-spline) regression, but assuming that the covariates are measured without error. Recently, a nested error regression model using a P-spline regression model, for the fixed part of the model, has been studied assuming measurement error in covariate in the Bayesian framework. In this paper, we propose a frequentist approach to study a semi-parametric nested error regression model using P-spline with a covariate measured with error. In particular, the pseudo-empirical best predictors of small-area means and their corresponding mean squared prediction error estimates are studied. Performance of the proposed approach is evaluated through a simulation and also by a real data application.

7 *Key words:* jackknife; linear mixed model; mean squared prediction error; penalised spline

8 1. Introduction

9 Sample surveys have been long used as a preferred means of gathering information about
10 a large population instead of census. Sometimes, to have a detailed analysis, estimating the
11 parameters for sub-populations within the overall population of interest is needed, but, due

* Author to whom correspondence should be addressed.

¹ Center for Healthcare Data Modeling, Departments of Biostatistics and Epidemiology, School of Public Health, Shahid Sadoughi University of Medical Sciences, Yazd, Iran

² Department of Biostatistics, Shahid Beheshti University of Medical Sciences, Tehran, Iran

³ Department of Community Health Sciences, University of Manitoba, Winnipeg, MB R3E 0W3, Canada

Email: Mahmoud.Torabi@umanitoba.ca

Acknowledgment. Constructive comments and suggestions of the Editor, an Associate Editor and two referees, which led to an improved version of this paper, are greatly appreciated. The research of Mahmoud Torabi was supported by the Natural Sciences and Engineering Research Council of Canada (NSERC) grant 368462-RGPIN.

12 to the cost and operational considerations, it is not always possible to have large enough
13 sample size to warrant accurate estimates for those sub-populations which are called small
14 areas. Since the traditional direct estimators do not provide adequate precision due to small
15 sample sizes corresponding to population sizes, the demand to use and develop small area
16 estimation methods has been greatly grown in recent years (Rao & Molina 2015). In order to
17 provide reliable estimates for areas with small or even zero sample sizes, "indirect" estimators
18 have been proposed in the context of small area estimation. The idea behind the indirect
19 estimators is to increase the effective sample size by borrowing strength from other sources
20 through a linking model using auxiliary information such as census data and administrative
21 data (Pfeffermann et al. 2013; Rao & Molina 2015; Jiang 2017).

22 Small area models, which are based on mixed model methodology, are divided into two broad
23 classes: (i) area-level models that relate the small area information on the response variable
24 to area-specific auxiliary variables, and (ii) unit-level models which relate the unit values of
25 the response variable to the unit-specific auxiliary variables with known area means and area-
26 specific covariates. Rao & Molina (2015) gave an extensive review of model-based small area
27 estimation under area-level and unit-level models. Focus of the current paper is on the unit-
28 level models.

29 One of the basic assumptions in unit-level models is that the covariates are measured
30 without error while this assumption may not be held in many real applications. Ignoring
31 the measurement error (ME) may cause the small area predictors perform worse than direct
32 estimators (Ybarra & Lohr 2008). In the context of classical ME model, there are two types
33 of ME models, functional and structural ME model. In the functional type, the unknown
34 true values of the covariate with ME are considered to be fixed which is in contrast with
35 the structural type where the unobserved covariate is assumed to be stochastic. Ghosh &
36 Sinha (2007), Datta, Rao & Torabi (2010), Torabi (2011) and Torkashvand, Jafari Jozani &
37 Torabi (2015) studied the functional ME for an area-specific covariate in the nested error
38 linear regression model. In these papers, the aim was to predict small area means with taking
39 into account functional ME in covariate. To estimate the ME in covariate, Ghosh & Sinha
40 (2007) proposed a moment estimator, Datta, Rao & Torabi (2010) suggested a maximum
41 likelihood estimator (MLE) and Torkashvand, Jafari Jozani & Torabi (2015) used a James-
42 Stein estimator to obtain pseudo-empirical Bayes (PEB) predictors of small area means.
43 Another basic assumption in the unit level model is that the mean of the continuous outcome
44 variable depends on the covariate value in a linear manner, while it might not hold in practice
45 and due to complexity of the relationship, assuming a linear trend might only be a crude
46 approximation. In such circumstances, parametric approaches will not properly work and to
47 express this relation, a semi-parametric smoothing method such as penalized spline (P-spline)

48 regression may be a good alternative (Eilers & Marx 1996). To see further applications of P-
49 spline, we refer readers to the overview of P-spline models written by Ruppert, Wand &
50 Carroll (2003).

51 In the context of small area estimation, Opsomer et al. (2008) studied P-spline regression
52 model in the linear mixed model set-up. Torabi & Shokoohi (2015) extended Opsomer et al.
53 (2008) model to a generalized linear mixed model (GLMM) to study normal and non-normal
54 responses. Shokoohi & Torabi (2018) studied the P-spline regression model in the class
55 of GLMMs to handle both time-series and cross-sectional response. Besides the P-spline
56 model, the non-parametric M-quantile regression has been also studied to model the non-
57 linear relationship between the q th M-quantile and the covariates in small area estimation
58 (Pratesi, Ranalli & Salvati 2008, 2009; Salvati, Ranalli & Pratesi 2011). Jiang, Nguyen & Rao
59 (2010) also proposed a procedure to select the small area model from a class of approximating
60 splines, using a fence method.

61 In practical applications, however, there are many situations in which not only the predictor
62 variable is not measured without error, but also the relationship between the response and
63 the covariate is not linear or it is even hard to find the relationship between the response
64 variable and the covariate. To deal with this problem, Hwang & Kim (2010, 2015) introduced
65 a non-parametric nested error regression model with truncated polynomial basis functions
66 and radial basis functions under functional ME model and predicted the small area means
67 via a Bayesian approach. Hwang & Kim (2016) extended their non-parametric model by
68 accommodating the covariates with and without ME again in a Bayesian framework.

69 In this paper, our aim is to take into account the functional ME in covariate in a semi-
70 parametric nested error regression model from a frequentist perspective. To that end, in
71 Section 2, we first rigorously study the model and present the "best" predictor of small area
72 means, which is the best linear unbiased predictor. We then estimate the true covariate, using
73 the maximum likelihood (ML) approach, to obtain the pseudo-"best" (PB) predictor of small
74 area means. We also obtain mean squared prediction error (MSPE) of PB predictor of small
75 area means. Furthermore, we use method-of-moments to estimate the model parameters to
76 derive pseudo-empirical "best" (PEB) predictor of small area means. To estimate the MSPE
77 of PEB predictor of small area means, we use the jackknife method. In order to evaluate our
78 proposed PEB predictor and its corresponding jackknife MSPE estimator, a simulation study
79 is conducted in Section 3. In Section 4, we employ the proposed model to predict the domain
80 (area) mean blood pressure measured in National Health and Nutrition Examination Survey
81 (NHANES), based on the cholesterol measured with error for some predefined domains,
82 which is an important national source of information examining the health status of the
83 population of the United States. Finally, we provide some concluding remarks in Section 5.

84

2. Model description

85

86

87

The nested error model with P-spline regression can be described as follows. Let y_{ij} be the variable of interest for the j -th unit ($j = 1, \dots, N_i$) at the i -th small area ($i = 1, \dots, m$) with corresponding observed covariate w_{ij} as

$$y_{ij} = f_0(x_i) + \nu_i + e_{ij} \quad (i = 1, \dots, m; j = 1, \dots, N_i), \quad (1)$$

$$w_{ij} = x_i + \eta_{ij} \quad (i = 1, \dots, m; j = 1, \dots, N_i), \quad (2)$$

88

89

90

91

92

where N_i is the population size of i -th area, x_i is a continuous covariate which is fixed but unknown, $\boldsymbol{\nu} = (\nu_1, \dots, \nu_m)^\top$ is the area-level random effects with $\nu_i \stackrel{iid}{\sim} N(0, \sigma_\nu^2)$, e_{ij} is the random error with $e_{ij} \stackrel{iid}{\sim} N(0, \sigma_e^2)$, η_{ij} is ME with $\eta_{ij} \stackrel{iid}{\sim} N(0, \sigma_\eta^2)$, and the function $f_0(x_i)$ is generally unknown. Note that in the context of ME, we do not observe x_i , but rather we observe w_{ij} as in (2). We can approximate $f_0(x_i)$ sufficiently well using P-spline as

$$f_0(x_i) = b_0 + b_1 x_i + \dots + b_p x_i^p + \sum_{a=1}^k \gamma_a (x_i - \tau_a)_+^p, \quad (3)$$

93

94

95

96

97

98

99

100

101

102

103

104

where p is the degree of spline, $(x)_+^p$ denotes the function $x^p I_{\{x>0\}}$, with I as the indicator function, $\{\tau_1, \dots, \tau_k\}$ is a set of knots which ties a sequence of line segments to trace the continuous relation between the covariate and the response variable, $\mathbf{b} = (b_0, \dots, b_p)^\top$ and $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_k)^\top$ are the regression coefficients of parameters and P-spline parts of the model, respectively. It is assumed that $\gamma_a \stackrel{iid}{\sim} N(0, \sigma_\gamma^2)$ and ν_i, e_{ij}, η_{ij} and γ_a are assumed to be mutually independent. Considering k large enough and defining the knots in a way that they vastly spread out over the range of x_i , this class of approximation is very comprehensive and can approximate most smooth functions. In this study, we determine the number of spline knots (k) as the minimum of 40 and number of x_i 's divided by 4, and the knots are quantiles of the distribution of x_i that are equally spaced (Ruppert 2002).

The goal is to predict the means of the response variable for the small areas of interest that is given by

$$\theta_i = \frac{1}{N_i} \sum_{j=1}^{N_i} y_{ij} \quad (i = 1, \dots, m),$$

105

106

107

108

on the basis of the sample data which are denoted by $\{(y_{ij}, w_{ij}); j = 1, \dots, n_i; i = 1, \dots, m\}$, where n_i is the sample size of i -th small area. It is assumed that the models (1) and (2) hold for the sample data assuming that there is no sample selection bias and the sampling design is not informative.

109 Clearly, as each of the P-spline and the small area models are in the class of random-effects
 110 models, the combination of these two models can also be treated as a linear mixed effects
 111 model as

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{Z}\boldsymbol{\gamma} + \mathbf{D}\boldsymbol{\nu} + \mathbf{e},$$

112 where $\mathbf{y} = (\mathbf{y}_1^\top, \dots, \mathbf{y}_m^\top)^\top$, $\mathbf{y}_i = (y_{i1}, \dots, y_{in_i})^\top$, and $n_T = \sum_{i=1}^m n_i$ is the
 113 total sample size. We define $\mathbf{X} = (\mathbf{x}_{11}, \dots, \mathbf{x}_{1n_1}, \dots, \mathbf{x}_{m1}, \dots, \mathbf{x}_{mn_m})^\top$ and
 114 $\mathbf{Z} = (\mathbf{z}_{11}, \dots, \mathbf{z}_{1n_1}, \dots, \mathbf{z}_{m1}, \dots, \mathbf{z}_{mn_m})^\top$ where $\mathbf{x}_{ij} = \mathbf{x}_i = (1, x_i, \dots, x_i^p)^\top$
 115 and $\mathbf{z}_{ij} = \mathbf{z}_i = ((x_i - \tau_1)_+^p, \dots, (x_i - \tau_k)_+^p)^\top$ are the vectors of the covariates
 116 for each sample in i -th area, respectively, for $j = 1, \dots, n_i$. We also define
 117 $\mathbf{D} = (\mathbf{d}_{11}, \dots, \mathbf{d}_{1n_1}, \dots, \mathbf{d}_{m1}, \dots, \mathbf{d}_{mn_m})^\top$ where $\mathbf{d}_{ij} = \mathbf{d}_i = (0, \dots, 0, 1, 0, \dots, 0)^\top$
 118 which the i -th element is equal to 1 and $\mathbf{e} = (\mathbf{e}_1^\top, \dots, \mathbf{e}_m^\top)^\top$, $\mathbf{e}_i = (e_{i1}, \dots, e_{in_i})^\top$.

119 2.1. Best predictor

120 To predict the i -th small area mean, first, assume that the value of the covariate x_i is not
 121 subject to the ME. Then, using the observed response data, the best predictor is given by

$$\begin{aligned} \hat{\theta}_i^B(x_i, \phi_1) &= N_i^{-1} \left[\sum_{j=1}^{n_i} y_{ij} + \sum_{j=n_i+1}^{N_i} \hat{y}_{ij} \right] \\ &= (1 - f_i) \bar{y}_i + f_i \left(b_0 + b_1 x_i + \dots + b_p x_i^p + \sum_{a=1}^k \hat{\gamma}_a (x_i - \tau_a)_+^p + \hat{\nu}_i \right), \end{aligned}$$

122 where $f_i = 1 - n_i/N_i$ for $i = 1, \dots, m$, $\hat{\boldsymbol{\gamma}} = \boldsymbol{\Sigma}_\gamma \mathbf{Z}^\top \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X}\mathbf{b})$, $\hat{\boldsymbol{\nu}} = \boldsymbol{\Sigma}_\nu \mathbf{D}^\top \mathbf{V}^{-1} (\mathbf{y} -$
 123 $\mathbf{X}\mathbf{b})$, which are the BLUP of the random effects $\boldsymbol{\gamma}$ and $\boldsymbol{\nu}$, respectively, where $\mathbf{V} =$
 124 $\text{var}(\mathbf{y}) = \mathbf{Z}\boldsymbol{\Sigma}_\gamma \mathbf{Z}^\top + \mathbf{D}\boldsymbol{\Sigma}_\nu \mathbf{D}^\top + \boldsymbol{\Sigma}_e$ with $\boldsymbol{\Sigma}_\gamma = \sigma_\gamma^2 \mathbf{I}_k$, $\boldsymbol{\Sigma}_\nu = \sigma_\nu^2 \mathbf{I}_m$, and $\boldsymbol{\Sigma}_e = \sigma_e^2 \mathbf{I}_{n_T}$ and
 125 $\phi_1 = (\mathbf{b}^\top, \sigma_e^2, \sigma_\nu^2, \sigma_\gamma^2)$ is assumed to be known. The corresponding *MSPE* of $\hat{\theta}_i^B$ is then
 126 given by

$$E(\hat{\theta}_i^B - \theta_i)^2 = f_i^2 \mathbf{q}_i^\top (\boldsymbol{\Sigma}_s - \boldsymbol{\Sigma}_s \boldsymbol{\Omega} \mathbf{V}^{-1} \boldsymbol{\Omega}^\top \boldsymbol{\Sigma}_s) \mathbf{q}_i,$$

127 where $\boldsymbol{\Omega} = (\mathbf{Z}, \mathbf{D})^\top$, $\mathbf{s} = (\boldsymbol{\gamma}^\top, \boldsymbol{\nu}^\top)^\top$, $\mathbf{q}_i = (\mathbf{z}_i^\top, \mathbf{l}_i^\top)^\top$, where \mathbf{l}_i is a vector with one as the
 128 i -th element and zero in other $(m - 1)$ elements, and $\boldsymbol{\Sigma}_s = \begin{bmatrix} \boldsymbol{\Sigma}_\gamma & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_\nu \end{bmatrix}$.

129 As it is clear, the introduced best predictor depends on x_i which may be measured with error
 130 in practice. As a result, the true values are not observed in such occasions. In the case of nested
 131 error linear regression model, Ghosh & Sinha (2007) estimated x_i with its moment estimator,
 132 \bar{w}_i . For the same set-up, Datta, Rao & Torabi (2010) proposed a MLE for x_i using all the

133 available data $\{y_{ij}, w_{ij}; j = 1, \dots, n_i; i = 1, \dots, m\}$ which is more efficient than Ghosh &
 134 Sinha (2007) estimator in terms of $MSPE$ of small area predictors. Following Datta, Rao &
 135 Torabi (2010), we use the MLE method to estimate the true value of the covariate x_i .

136 2.2. Pseudo-best predictor

137 In order to estimate x_i , we consider the sample means of our data $\{(y_{ij}, w_{ij}); j =$
 138 $1, \dots, n_i; i = 1, \dots, m\}$ through the equations (1)-(3) as

$$\begin{aligned} \bar{y}_i &= b_0 + b_1 x_i + \dots + b_p x_i^p + \sum_{a=1}^k \gamma_a (x_i - \tau_a)_+^p + \nu_i + \bar{e}_i \\ &= b_0 + b_1 x_i + \dots + b_p x_i^p + \bar{u}_{1i}, \end{aligned} \quad (4)$$

$$\bar{w}_i = x_i + \bar{u}_{2i}, \quad (5)$$

139 where $\bar{u}_{1i} \sim N\left(0, \sigma_\gamma^2 \sum_{a=1}^k (x_i - \tau_a)_+^{2p} + \sigma_\nu^2 + \frac{\sigma_e^2}{n_i}\right)$ and $\bar{u}_{2i} = \bar{\eta}_i \sim N\left(0, \frac{\sigma_\eta^2}{n_i}\right)$. Since \bar{u}_{1i}
 140 is independent of \bar{u}_{2i} , the log-likelihood function, $l(x_i)$, can be expressed as the log of joint
 141 density $f(\bar{y}_i, \bar{w}_i | x_i) = f(\bar{y}_i | x_i) f(\bar{w}_i | x_i)$ through

$$\begin{aligned} l(x_i) &= \log(f(\bar{y}_i, \bar{w}_i | x_i)) \\ &\propto -\frac{1}{2\sigma_\eta^2} n_i (\bar{w}_i - x_i)^2 - \frac{1}{2} \log \left[2\pi \left(\sigma_\gamma^2 \sum_{a=1}^k (x_i - \tau_a)_+^{2p} + \sigma_\nu^2 + \frac{\sigma_e^2}{n_i} \right) \right] \\ &\quad - \frac{(\bar{y}_i - b_0 - b_1 x_i - \dots - b_p x_i^p)^2}{2 \left[\sigma_\gamma^2 \sum_{a=1}^k (x_i - \tau_a)_+^{2p} + \sigma_\nu^2 + \frac{\sigma_e^2}{n_i} \right]}. \end{aligned}$$

142 Since $l(x_i)$ does not have a closed form, we use numerical methods for maximization.
 143 Maximizing the likelihood function with respect to x_i and substituting \tilde{x}_i , which is the
 144 estimate of x_i , for x_i in the best estimator leads to the following PB predictor

$$\begin{aligned} \hat{\theta}_i^{PB} &= \hat{\theta}_i^{PB}(\phi) \\ &= (1 - f_i) \bar{y}_i + f_i (b_0 + b_1 \tilde{x}_i + \dots + b_p \tilde{x}_i^p + \sum_{a=1}^k \hat{\gamma}_a (\tilde{x}_i - \tau_a)_+^p + \hat{\nu}_i), \end{aligned}$$

145 where $\phi = (\phi_1, \sigma_\eta^2)$. Since $E(\hat{\theta}_i^{PB} - \theta_i | \mathbf{y}_i) = 0$, the $MSPE$ of PB predictor is given by

$$\begin{aligned} MSPE(\hat{\theta}_i^{PB}) &= E(\hat{\theta}_i^{PB} - \theta_i)^2 \\ &= f_i^2 E \left\{ b_1 (\tilde{x}_i - x_i) + \dots + b_p (\tilde{x}_i^p - x_i^p) + \sum_{a=1}^k \hat{\gamma}_a \left[(\tilde{x}_i - \tau_a)_+^p - (x_i - \tau_a)_+^p \right] \right\}^2 \\ &\quad + f_i^2 \mathbf{q}_i^\top (\boldsymbol{\Sigma}_s - \boldsymbol{\Sigma}_s \boldsymbol{\Omega} \mathbf{V}^{-1} \boldsymbol{\Omega}^\top \boldsymbol{\Sigma}_s) \mathbf{q}_i \equiv g_{1i}(\phi). \end{aligned}$$

146 In reality, the PB predictor and corresponding $MSPPE$ are not computable as they depend
147 on the model parameters ϕ .

148 2.3. Pseudo-empirical best predictor

149 In order to predict the small area means we now need to estimate ϕ . Then substituting $\hat{\phi}$
150 for ϕ in $\hat{\theta}_i^{PEB}$ gives PEB predictor, $\hat{\theta}_i^{PEB}$, of small area means. Here, we use the method-of-
151 moments to estimate the model parameters ϕ . Following Ghosh & Sinha (2007), the estimates
152 of random error and ME variances are

$$\hat{\sigma}_e^2 = \frac{SSW_{\mathbf{y}}}{n_T - m},$$

$$\hat{\sigma}_\eta^2 = \frac{SSW_{\mathbf{w}}}{n_T - m},$$

153 where

$$SSW_{\mathbf{y}} = \sum_{i=1}^m \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2,$$

$$SSW_{\mathbf{w}} = \sum_{i=1}^m \sum_{j=1}^{n_i} (w_{ij} - \bar{w}_i)^2.$$

154 In the next step, we need to estimate the regression coefficients \mathbf{b} . Based on equations (4)-(5),
155 we can write the model as

$$\bar{y}_i = \bar{\mathbf{x}}_i^\top \mathbf{b} + \bar{\mathbf{z}}_i^\top \boldsymbol{\gamma} + \bar{\mathbf{d}}_i^\top \boldsymbol{\nu} + \bar{e}_i, \quad (6)$$

$$\mathbf{h}_i = \bar{\mathbf{x}}_i + \bar{\boldsymbol{\eta}}_i, \quad (7)$$

156 where $\mathbf{h}_i = (1, \bar{w}_i, \dots, \bar{w}_i^p)^\top$ and $\bar{\boldsymbol{\eta}}_i = (\bar{\eta}_{0i}, \bar{\eta}_{1i}, \dots, \bar{\eta}_{pi})^\top$. Since $E(\bar{e}_i) = 0$, $E(\bar{\boldsymbol{\eta}}_i) = \mathbf{0}$,
157 $var(\bar{e}_i) = \sigma_e^2/n_i$, and $var(\bar{\boldsymbol{\eta}}_i) = n_i^{-1}\boldsymbol{\Sigma}_\eta$ with $\boldsymbol{\Sigma}_\eta$ as the variance-covariance matrix of
158 measurement errors, we have

$$E(\mathbf{h}_i) = \bar{\mathbf{x}}_i,$$

159

$$E(n_i \mathbf{h}_i \mathbf{h}_i^\top) = n_i \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top + \boldsymbol{\Sigma}_\eta$$

160 From equation (6), we also have

$$n_i \bar{\mathbf{x}}_i \bar{y}_i = n_i \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \mathbf{b} + n_i \bar{\mathbf{x}}_i \bar{\mathbf{z}}_i^\top \boldsymbol{\gamma} + n_i \bar{\mathbf{x}}_i \bar{\mathbf{d}}_i^\top \boldsymbol{\nu} + n_i \bar{\mathbf{x}}_i \bar{e}_i \quad (8)$$

161 and taking expectation and then taking expectation from both sides of the equation (8) leads
 162 to

$$\begin{aligned} n_i \bar{\mathbf{x}}_i E(\bar{y}_i) &= n_i \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \mathbf{b}, \\ \frac{1}{m} \sum_{i=1}^m n_i \bar{\mathbf{x}}_i E(\bar{y}_i) &= \left(\frac{1}{m} \sum_{i=1}^m n_i \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \right) \mathbf{b}, \\ \mathbf{b} &= \left(\frac{1}{m} \sum_{i=1}^m n_i \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \right)^{-1} \left(\frac{1}{m} \sum_{i=1}^m n_i \bar{\mathbf{x}}_i E(\bar{y}_i) \right). \end{aligned}$$

163 The plug-in estimator of \mathbf{b} is then given by

$$\hat{\mathbf{b}} = \left(\frac{1}{m} \sum_{i=1}^m (n_i \mathbf{h}_i \mathbf{h}_i^\top - \hat{\Sigma}_\eta) \right)^{-1} \left(\frac{1}{m} \sum_{i=1}^m n_i \mathbf{h}_i \bar{y}_i \right).$$

164 To estimate σ_γ^2 and σ_ν^2 , let $MSB_{\mathbf{y}} = (m-1)^{-1} \sum_{i=1}^m n_i (\bar{y}_i - \bar{y})^2$, then, following Datta,
 165 Rao & Torabi (2010), we have

$$\begin{aligned} E(MSB_{\mathbf{y}}) &= \frac{\sigma_\gamma^2}{m-1} \left[\sum_{i=1}^m \left(n_i - \frac{n_i^2}{n_T} \right) \sum_{a=1}^k (x_i - \tau_a)_+^{2p} \right] + \frac{\sigma_\nu^2}{m-1} \sum_{i=1}^m \left(n_i - \frac{n_i^2}{n_T} \right) + \sigma_e^2 + \\ &\quad \frac{1}{m-1} \sum_{i=1}^m n_i \left[b_1(x_i - \bar{x}) + \dots + b_p(x_i^p - \bar{x}^p) \right]^2, \end{aligned}$$

166 and the estimates of σ_γ^2 and σ_ν^2 are then obtained as

$$\begin{aligned} \hat{\sigma}_\gamma^2 &= \left\{ (m-1)MSB_{\mathbf{y}} - (m-1)\hat{\sigma}_e^2 - \sum_{i=1}^m n_i \left[\hat{b}_1(\hat{x}_i - \bar{x}) + \dots + \hat{b}_p(\hat{x}_i^p - \bar{x}^p) \right]^2 \right. \\ &\quad \left. - \hat{\sigma}_\nu^2 \sum_{i=1}^m \left(n_i - \frac{n_i^2}{n_T} \right) \right\} / \left\{ \sum_{i=1}^m \left(n_i - \frac{n_i^2}{n_T} \right) \sum_{a=1}^k (\hat{x}_i - \tau_a)_+^{2p} \right\}, \end{aligned}$$

167 and

$$\begin{aligned} \hat{\sigma}_\nu^2 &= \frac{1}{\sum_{i=1}^m \left(n_i - \frac{n_i^2}{n_T} \right)} \left\{ (m-1)MSB_{\mathbf{y}} - (m-1)\hat{\sigma}_e^2 - \sum_{i=1}^m n_i \left[\hat{b}_1(\hat{x}_i - \bar{x}) + \dots + \hat{b}_p(\hat{x}_i^p - \bar{x}^p) \right]^2 \right. \\ &\quad \left. - \hat{\sigma}_\gamma^2 \left[\sum_{i=1}^m \left(n_i - \frac{n_i^2}{n_T} \right) \sum_{a=1}^k (\hat{x}_i - \tau_a)_+^{2p} \right] \right\}, \end{aligned}$$

168 where $\bar{x}^j = \frac{1}{n_T} \sum_{i=1}^m n_i \hat{x}_i^j$, ($j = 1, \dots, p$). Since the estimates of $(x_1, \dots, x_p, \sigma_\nu^2, \sigma_\gamma^2)$ are
 169 dependent of each other, they are estimated via an iterative algorithm. In particular, to
 170 estimate (x_1, \dots, x_p) the above two equations related to $\hat{\sigma}_\gamma^2$ and $\hat{\sigma}_\nu^2$ are estimated iteratively
 171 with respect to (x_1, \dots, x_p) through the within mean square error of the observed values of
 172 the covariate. Finally, the *PEB* predictor of θ_i is given by

$$\begin{aligned} \hat{\theta}_i^{PEB} &= \hat{\theta}_i^{PB}(\hat{\phi}) \\ &= (1 - f_i)\bar{y}_i + f_i \left(b_0 + b_1\tilde{x}_i + \dots + b_p\tilde{x}_i^p + \sum_{a=1}^k \hat{\gamma}_a(\tilde{x}_i - \tau_a)_+^p + \hat{\nu}_i \right). \end{aligned}$$

173 To measure the variability of $\hat{\theta}_i^{PEB}$, the *MSPE* of *PEB* predictor can be decomposed as

$$\begin{aligned} MSPE(\hat{\theta}_i^{PEB}) &= E(\hat{\theta}_i^{PB} - \theta_i)^2 + E(\hat{\theta}_i^{PEB} - \hat{\theta}_i^{PB})^2 + 2E\left[(\hat{\theta}_i^{PB} - \theta_i)(\hat{\theta}_i^{PEB} - \hat{\theta}_i^{PB})\right] \\ &= M_{1i} + M_{2i} + 2M_{3i}, \end{aligned}$$

174 where $M_{1i} = g_{1i}(\phi)$. The *MSPE* of *PEB* predictor is not computable as it depends on
 175 model parameters. In order to estimate the *MSPE* of *PEB* predictor of small area means,
 176 we apply the jackknife method which was proposed by Jiang et al. (2002) and Chen & Lahiri
 177 (2002). To be able to use the jackknife method, similar to other studies which were done in the
 178 context of functional ME for area-level and unit-level models such as Ybarra & Lohr (2008),
 179 Datta, Rao & Torabi (2010) and Torkashvand, Jafari Jozani & Torabi (2015), we approximate
 180 the *MSPE* as

$$MSPE(\hat{\theta}_i^{PEB}) \approx M_{1i} + M_{2i},$$

181 by ignoring the cross-product term, noting that there is no-closed form expression for the
 182 $MSPE(\hat{\theta}_i^{PEB})$ (Haslett & Welsh 2019). We will report the magnitude of M_{3i} in the
 183 simulation study section.

184 To estimate M_{1i} , a jackknife bias correction is used which is given by

$$\hat{M}_{1iJ} = g_{1i}(\hat{\phi}) - \sum_{l \neq i} \psi_l \left[g_{1i}(\hat{\phi}_{-l}) - g_{1i}(\hat{\phi}) \right], \quad l = 1, \dots, m, \quad (9)$$

185 where $\psi_l = 1 + O(m^{-1})$ is a suitable weight (Chen & Lahiri 2002). Here, $g_{1i}(\hat{\phi})$ is the plug-
 186 in estimator of $g_{1i}(\phi)$ and $\hat{\phi}_{-l}$ is the moment estimator of ϕ , obtained by omitting the l -th
 187 area data set from the full data set $\{(y_{ij}, w_{ij}); j = 1, \dots, n_i; i = 1, \dots, m\}$. This is done for
 188 each $l \neq i$ (except the i -th area) to get $m - 1$ estimators for ϕ .

189 The jackknife estimator of M_{2i} is given by

$$\hat{M}_{2iJ} = \sum_{l \neq i} \psi_l (\hat{\theta}_{i,-l}^{PEB} - \hat{\theta}_i^{PEB})^2, \quad l = 1, \dots, m, \quad (10)$$

190 where $\hat{\theta}_{i,-l}^{PEB}$ is the plug-in estimator of $\hat{\theta}_i^{PEB}$, in which the vector of the parameters (ϕ)
 191 is estimated by deleting the l -th area data set from the full data set each time. Finally, the
 192 jackknife estimator of $MSPE(\hat{\theta}_i^{PEB})$ is obtained by taking the sum of (9) and (10) as

$$mspe_J(\hat{\theta}_i^{PEB}) = \hat{M}_{1iJ} + \hat{M}_{2iJ}.$$

193 Assuming $\psi_l = 1 - \mathbf{h}_l^\top (\sum_{t \neq i} \mathbf{h}_t \mathbf{h}_t^\top)^{-1} \mathbf{h}_l$ and $\psi_l = \frac{m-2}{m-1}$, the weighted and unweighted
 194 versions of jackknife estimator of $MSPE(\hat{\theta}_i^{PEB})$ are obtained, respectively. Note that
 195 in small area estimation (Rao & Molina 2015), the notation *mspe* is usually used as the
 196 estimator of $MSPE$.

197

198 3. Simulation study

199 In this section, we carry out a simulation study to compare the performance of the
 200 proposed approach in the P-spline model which takes into account the ME in the area level
 201 predictor variable and the P-spline model which ignores the ME (naive model). To this end,
 202 the population responses are generated from the model (1) with three choices for $f_0(x_i)$:
 203 linear, quadratic, and exponential 1 (see later of this section for details of choices of $f_0(x_i)$).
 204 Note that in small area estimation, the asymptotic result will apply for large number of small
 205 areas m . So, for large m the effects of model parameter estimate would be vanished as long
 206 as the model parameter estimators are consistent. Hence, it is important in the context of
 207 small area estimation to show how good is the proposed model for finite sample (small m).
 208 Therefore, we assume that the population units are distributed across $m = 40$ areas equally
 209 in a way that $N_i = 400$, ($i = 1, \dots, m$), and equal sample sizes are taken from each area
 210 as $n_i = 4$, ($i = 1, \dots, m$). We generate $R = 1000$ independent sets of $\{\nu_i^{(r)}; i = 1, \dots, m\}$,
 211 $\{e_{ij}^{(r)}; j = 1, \dots, N_i; i = 1, \dots, m\}$ from Normal distribution with mean zero and variance
 212 σ_ν^2 and σ_e^2 , respectively. We assume $\sigma_\nu^2 = 1$ and $\sigma_e^2 = 1$ for linear case and $\sigma_\nu^2 = 0.1$ and
 213 $\sigma_e^2 = 0.3$ for non-linear cases. The true values of the predictor $\{x_i; i = 1, \dots, m\}$ are also
 214 generated from a uniform distribution between 10 and 30 for linear case and uniform
 215 distribution between -3 and 3 for non-linear cases and treat them fixed through the simulation
 216 study. Using $\{x_i, \nu_i^{(r)}, e_{ij}^{(r)}\}$, the population responses are generated from the model (1) with

217 three choices of $f_0(x_i)$ (linear, quadratic, and exponential 1) as

$$\begin{aligned} y_{ij}^{(r)} &= 1 + x_i + \nu_i^{(r)} + e_{ij}^{(r)}, \quad r = 1, \dots, R; j = 1, \dots, N_i; i = 1, \dots, m, \\ y_{ij}^{(r)} &= 0.4 + 0.4x_i - 0.65x_i^2 + \nu_i^{(r)} + e_{ij}^{(r)}, \\ y_{ij}^{(r)} &= 1 + x_i - 0.7 \exp(x_i) + \nu_i^{(r)} + e_{ij}^{(r)}, \end{aligned}$$

218 following Breidt, Claeskens & Opsomer (2005), Rao, Sinha & Dumitrescu (2014) and
219 Shokoohi & Torabi (2018). The population mean response of i -th area for r -th simulation
220 is given by

$$\theta_i^{(r)} = \frac{1}{N_i} \sum_{j=1}^{N_i} y_{ij}^{(r)}.$$

221 We then generate simple random samples from each simulated population responses.
222 Furthermore, we generate the observed values of the predictor variable from the ME model
223 $w_{ij}^{(r)} = x_i + \eta_{ij}^{(r)}$, ($j = 1, \dots, n_i; i = 1, \dots, m$), where $\eta_{ij}^{(r)}$ is generated from a normal
224 distribution with mean zero and variance $\sigma_\eta^2 = 2$ for linear case and $\sigma_\eta^2 = 0.6$ for non-linear
225 quadratic and exponential 1 cases. Thereafter, for each simulated data set $\{(w_{ij}^{(r)}, y_{ij}^{(r)}); j =$
226 $1, \dots, n_i; i = 1, \dots, m\}$ in each scenario (linear, quadratic, and exponential 1), the P-spline
227 model with and without ME are fitted assuming $p = 1$ (which has piecewise linear fit, with
228 the changes in slope at each knot regarded as random with variance σ_γ^2). Table 1 presents
229 the moment estimators of the proposed model parameters for each scenario. In particular, in
230 the case of linear model as the true model, the model provides the regression coefficients
231 which are very close to the true values; note that the variance component of P-spline (σ_γ^2) is
232 close to zero since the true model is linear. Also, in the cases of quadratic and exponential
233 1 models, the fitted model works very well to track the true values, note that in these two
234 scenarios, we only need to compare the variance components estimates of the model with the
235 corresponding true values. We observe that the variations of the models from the linearity
236 (quadratic and exponential 1) are well captured through the estimate of σ_γ^2 in the proposed
237 model.

238 [Table 1 about here]

239 Figures 1a, 1b, and 1c show PEB predictions of small area means. It seems that both models
240 have similar predictions for small area means.

241 [Figure 1 about here]

242 Now we need to evaluate the accuracy of $\hat{\theta}_i^{PEB}$ in our proposed approach. To that end,
 243 we calculate the empirical $MSPE$ ($EMSPPE$) of $\hat{\theta}_i^{PEB}$ which is given by

$$EMSPPE(\hat{\theta}_i^{PEB(r)}) = \frac{1}{R} \sum_{r=1}^R (\hat{\theta}_i^{PEB(r)} - \theta_i^{(r)})^2.$$

244 To also evaluate the magnitude of the cross-product term involved in the $MSPE$ of PEB
 245 predictor of small area means, $MSPE$ of $\hat{\theta}_i^{PEB(r)}$ can be decomposed as

$$EMSPPE(\hat{\theta}_i^{PEB(r)}) = M_{ip} = M_{1ip} + M_{2ip} + 2M_{3ip},$$

246 where $M_{1ip} = R^{-1} \sum_{r=1}^R (\hat{\theta}_i^{PB(r)} - \theta_i^{(r)})^2$, $M_{2ip} = R^{-1} \sum_{r=1}^R (\hat{\theta}_i^{PEB(r)} - \hat{\theta}_i^{PB(r)})^2$ and
 247 $M_{3ip} = R^{-1} \sum_{r=1}^R (\hat{\theta}_i^{PB(r)} - \theta_i^{(r)})(\hat{\theta}_i^{PEB(r)} - \hat{\theta}_i^{PB(r)})$.

248 [Figure 2 about here]

249 [Figure 3 about here]

250 [Figure 4 about here]

251 Figures 2a, 3a, and 4a show the $EMSPPE$ of $\hat{\theta}_i^{PEB(r)}$ and its decomposition for the both
 252 proposed and naive models for three cases (linear, quadratic, and exponential 1). Based on the
 253 results for the proposed and naive models for all three cases (linear, quadratic, and exponential
 254 1), it seems that the PEB predictions of small area means are near the true means since the
 255 values of M_{ip} are close to zero particularly in the cases of quadratic and exponential forms,
 256 and the most contribution of $MSPE$ is attributed to M_{1ip} as expected. It is also clear from
 257 Figures 2a, 3a, and 4a that the contribution of the cross-product M_{3ip} term involved in the
 258 $MSPE$ is small in the proposed model compared to the naive model.

259 In order to evaluate the performance of the weighted, $m_{spe_{JW}}(\hat{\theta}_i^{PEB})$, and unweighted,
 260 $m_{spe_J}(\hat{\theta}_i^{PEB})$, jackknife $MSPE$ estimators, the empirical relative bias (RB) of these
 261 estimators are computed through the following expression

$$RB_i = \frac{E(m_{spe_i})}{EMSPPE_i} - 1 \quad (i = 1, \dots, m),$$

262 where $E(m_{spe_i})$ is the average of simulated jackknife $MSPE$ estimate of PEB predictor
 263 of small area mean i . The results of RB of jackknife m_{spe} of PEB predictor of small area
 264 means (weighted and unweighted) for the proposed and naive models are also reported for
 265 all three cases (linear, quadratic, and exponential 1) in Figures 2b, 3b, and 4b. As we expect,
 266 the proposed P-spline model performs very well in terms of RB for the three scenarios in
 267 this simulation set-up, and it appears that the unweighted $MSPE$ estimates perform as well

268 as the corresponding weighted version. Based on our empirical findings, jackknife method
269 causes serious overestimation of $MSPE$ of the PEB predictor in naive model due to the
270 large values of M_{3ip} as shown in Figures 2a to 4a. Furthermore, in order to evaluate the effect
271 of level of measurement error on estimates and predictions, we consider different values of
272 measurement error variance for linear and non-linear cases. Table 2 presents the moment
273 estimators of the proposed model parameters for each scenario. As shown in Table 2, our
274 proposed P-spline model works very well in terms of model parameters estimate for different
275 spline forms and measurement error variances.

276 [Table 2 about here]

277 In terms of performance of naive model, for example, we observed that in the case of linear
278 model, with increasing the value of measurement error variance, parameter estimate of slope
279 is attenuated (not shown here) unlike the proposed P-spline model. The bias in the slope
280 estimate which is caused by measurement error is discussed in the literature (Carroll 2006).
281 Figures 4 to 6 show $EMSPPE$ of PEB predictors and its components for the proposed and
282 naive models and related percent relative bias of jackknife estimators of unweighted and
283 weighted $MSPE$ for the proposed and naive models. With increasing measurement error
284 variance σ_{η}^2 , the proposed P-spline model still shows good performance in terms of RB in
285 all three scenarios. Note that in case of naive model with exponential form, the cross-product
286 term (M_{3in}) has the same magnitude as the leading term (M_{1in}) but with opposite sign. This
287 is the reason that RB for the naive model is also as good as the proposed model in the case
288 of exponential model.

289 [Figure 5 about here]

290 [Figure 6 about here]

291 [Figure 7 about here]

292 4. Application

293 In this section, we employ our proposed P-spline model to analyze data from the 2013–
294 2014 US NHANES. The NHANES is a yearly survey to determine the health and nutritional
295 status of adults and children in the United States. According to the literature, there is a
296 significant positive relationship between obesity and blood pressure (Lee, Bacha & Arslanian
297 2006; Choy et al. 2011; Duncan et al. 2013). Since waist circumference (WC) index is
298 expressed as the main indicator of abdominal fat accumulation, hence, in this study, our aim is
299 to predict the mean systolic blood pressure in some demographic domains of interest using the
300 WC of NHANES participants as an auxiliary information which is likely measured with error

301 (Caballero 2005). The focus of our analysis is on 5588 participants. We build fifty domains
 302 ($m = 50$) with sample sizes ranging from 31 to 479, based on sex, five age categories (20-
 303 29, 30-39, 40-48, 50-59, and 60-84), and five race and ethnicity groups (Mexican American,
 304 Other Hispanic, White non-Hispanic, Black non-Hispanic and Other), and use the WC and
 305 systolic blood pressure as the values of the predictor and response variables. Figure 8 shows
 306 the mean systolic blood pressure versus the mean WC in 50 domains. It appears from Figure 8
 307 that there is a non-linear relationship between these two variables. So, semi-parametric
 308 models such as P-spline models are good candidates to analyze this data set. In addition, since
 309 the WC is prone to measurement error, applying the proposed model (P-spline with $p = 1$
 310 which considers ME of the variable WC) seems to be worthwhile. To compare the proposed
 311 model with the naive model, we also analyze this data set with ignoring measurement error.
 312 The estimated parameters (and standard errors using jackknife method) for the proposed
 313 model are $\hat{b}_0 = 35.92(20.41)$, $\hat{b}_1 = 0.88(0.20)$, $\hat{\sigma}_\eta^2 = 234.33(13.07)$, $\hat{\sigma}_e^2 = 249.27(23.13)$,
 314 $\hat{\sigma}_\nu^2 = 5.00(6 \times 10^{-11})$, $\hat{\sigma}_\gamma^2 = 0.13(0.05)$ and for the naive model are $\hat{b}_0 = 42.74(18.53)$,
 315 $\hat{b}_1 = 0.81(0.18)$, $\hat{\sigma}_e^2 = 249.27(23.13)$, $\hat{\sigma}_\nu^2 = 5.00(3 \times 10^{-11})$, $\hat{\sigma}_\gamma^2 = 0.12(0.05)$. From the
 316 estimated parameter and corresponding standard error of variance of ME and comparing
 317 the test statistic with critical value 1.96, it is observed that the WC is measured with error.
 318 Furthermore, based on the obtained significant non-zero σ_γ^2 , it is clear that there is a non-
 319 linear relationship between WC and systolic blood pressure. So, it seems that neither the
 320 non-linear relation nor the measurement error in WC, which causes attenuation in estimate
 321 of the slope, can be ignored. It is also worth mentioning that the proposed model shows
 322 that the WC has a positive effect in predicting blood pressure which is supported by the
 323 literature (Lee, Bacha & Arslanian 2006; Choy et al. 2011; Duncan et al. 2013). The boxplots
 324 of *PEB* predictor of mean blood pressure for predefined domains, and their weighted and
 325 unweighted jackknife estimates of *MSPE* for both models are presented in Figures 9a and
 326 9b based on NHANES study. According to Figure 9a, it appears that both models behave
 327 similarly to predict mean blood pressure for predefined domains generally. We observe that
 328 mean blood pressure in men is higher than women fixing age, race and ethnicity. To study
 329 the effect of race and ethnicity, it is seen that mean blood pressure has approximately the
 330 same level in Mexican American, Other Hispanic, White non-Hispanic categories, while
 331 Black non-Hispanic and Other category have lower and White non-Hispanic group has higher
 332 level of blood pressure. Furthermore, based on the results, increasing age leads to a higher
 333 blood pressure, fixing the other two variables. In terms of *MSPE* estimation, the naive
 334 jackknife estimation of *MSPE* behave differently than the proposed model. In general, the
 335 weighted and unweighted estimators of *MSPE* for the proposed model are smaller than the
 336 corresponding naive estimators of *MSPE*. Based on the simulation results, we can conclude

337 that ignoring the ME may lead to wrong conclusions in terms of overestimation in jackknife
338 $MSPE$ estimation of small area mean predictors.

339 [Figure 8 about here]

340 [Figure 9 about here]

341 5. Discussion

342 We have proposed a semi-parametric nested error regression model with functional ME
343 in area-level covariate in a frequentist framework. According to the moment estimators in
344 simulation part, it is observed that in the case of linear model as the true model, the regression
345 coefficients of the proposed P-spline model are so close to the true values. Also, the estimate
346 of variance of spline term is near zero in linear case while this estimate is non-zero for non-
347 linear cases which indicates that the proposed model detects the existence of a non-linear
348 relationship between the response and predictor variables. In particular, we have derived
349 the PEB predictor of small area means and obtained the corresponding $MSPE$ of PEB
350 predictor of small area means. We have also proposed jackknife estimators of the $MSPE$ of
351 the PEB predictors. We have shown through a simulation that although the PEB predictor
352 of small area means are very similar for the both proposed and naive methods, however, our
353 proposed approach works very well in terms of jackknife $MSPE$ estimates of the PEB
354 predictor of small area means compared to the naive model which ignores the ME in the
355 covariate that causes serious overestimation due to the large value of cross-product terms
356 involved in the $MSPE$ of PEB predictor of small area means. We have also studied the effect
357 of increasing ME in predictor variable for both models in different scenarios with considering
358 different values for variance of ME. We have observed increasing ME causes attenuation in
359 slope estimate in naive model unlike the proposed model.

360 Our proposed model is developed based on one covariate. An extension of our work
361 to multiple covariates measured with error is simple, however, it will add much more
362 unnecessary complexity to the model as one needs to define different spline terms for each
363 covariate. One can also extend our approach to deal with non-normal random effects (Hui,
364 Muller & Welsh 2020). As an extension of the proposed model, survey weights (Torabi 2011)
365 can also be used in the estimation process in order to increase the efficiency of the PEB
366 predictor of small area means. In this paper, we have assumed that the size of sample in the
367 response and observed covariate is the same in each small area, however, one can extend our
368 proposed approach and use multiple source of data with different sample sizes (Datta et al.
369 2018). One can extend our proposed model to generalized linear models (Torabi & Shokoohi
370 2015). One can also extend our proposed model and use the bootstrap and simple, unified,

371 Monte-Carlo assisted (Sumca) methods as alternatives for $MSPE$ estimation of small area
372 mean predictors (Jiang & Torabi 2020). These are some of the topics for future study.

373

Appendix

374 The supplementary materials provide R codes and corresponding “readme” files for the
375 simulation and real application conducted in this paper.

References

376

- 377 BREIDT, F., CLAESKENS, G. & OPSOMER, J. (2005). Model-assisted estimation for complex surveys using
378 penalised splines. *Biometrika* **92**, 831–846.
- 379 CABALLERO, B. (2005). *Encyclopedia of Human Nutrition*. Oxford: Elsevier Academic Press, 2nd edn.
- 380 CARROLL, RAYMOND J., D.R.L.A.S.C.M.C. (2006). *Measurement error in nonlinear models: a modern*
381 *perspective*. Chapman and Hall/CRC, 2nd edn.
- 382 CHEN, S. & LAHIRI, P. (2002). On mean squared prediction error estimation in small area estimation
383 problems. In *Proceedings of the Survey Research Methods Section* , 473–477.
- 384 CHOY, C.S., CHAN, W.Y., CHEN, T.L., SHIH, C.C., WU, L.C. & LIAO, C.C. (2011). Waist circumference
385 and risk of elevated blood pressure in children: a cross-sectional study. *BMC Public Health* **11**, 1–7.
- 386 DATTA, G.S., RAO, J. & TORABI, M. (2010). Pseudo-empirical bayes estimation of small area means
387 under a nested error linear regression model with functional measurement errors. *Journal of Statistical*
388 *Planning and Inference* **140**, 2952–2962.
- 389 DATTA, G.S., TORABI, M., RAO, J. & LIU, B. (2018). Small area estimation with multiple covariates
390 measured with errors: A nested error linear regression approach of combining multiple surveys. *Journal*
391 *of Multivariate Analysis* **167**, 49–59.
- 392 DUNCAN, M.J., MOTA, J., VALE, S., SANTOS, M.P. & RIBEIRO, J.C. (2013). Associations between
393 body mass index, waist circumference and body shape index with resting blood pressure in portuguese
394 adolescents. *Annals of Human Biology* **40**, 163–167.
- 395 EILERS, P.H. & MARX, B.D. (1996). Flexible smoothing with b-splines and penalties. *Statistical science* ,
396 89–102.
- 397 GHOSH, M. & SINHA, K. (2007). Empirical bayes estimation in finite population sampling under functional
398 measurement error models. *Journal of statistical planning and inference* **137**, 2759–2773.
- 399 HASLETT, S. & WELSH, A. (2019). Eblup: Empirical best linear unbiased estimation. *Wiley StatsRef,*
400 *Statistics Reference Online* .
- 401 HUI, F., MULLER, S. & WELSH, A. (2020). Random effects misspecification can have severe consequences
402 for random effects inference in linear mixed models. *International Statistical Review* , 1–21.
- 403 HWANG, J. & KIM, D.H. (2015). Bayesian curve-fitting with radial basis functions under functional
404 measurement error model. *Journal of the Korean Data and Information Science Society* **26**, 749–754.
- 405 HWANG, J. & KIM, D.H. (2016). Multivariable bayesian curve-fitting under functional measurement error
406 model. *Journal of the Korean Data and Information Science Society* **27**, 1645–1651.
- 407 HWANG, J.S. & KIM, D.H. (2010). Semiparametric bayesian estimation under functional measurement error
408 model. *Journal of the Korean Data and Information Science Society* **21**, 379–385.
- 409 JIANG, J. (2017). *Asymptotic analysis of mixed effects models: theory, applications, and open problems*.
410 Chapman and Hall/CRC.
- 411 JIANG, J., LAHIRI, P., WAN, S.M. et al. (2002). A unified jackknife theory for empirical best prediction
412 with m-estimation. *The Annals of Statistics* **30**, 1782–1810.
- 413 JIANG, J., NGUYEN, T. & RAO, J.S. (2010). Fence method for nonparametric small area estimation. *Survey*
414 *Methodology* **36**, 3–11.
- 415 JIANG, J. & TORABI, M. (2020). Sumca: simple, unified, monte-carlo-assisted approach to second-order
416 unbiased mean-squared prediction error estimation. *Journal of the Royal Statistical Society: Series B*
417 *(Statistical Methodology)* .
- 418 LEE, S., BACHA, F. & ARSLANIAN, S.A. (2006). Waist circumference, blood pressure, and lipid
419 components of the metabolic syndrome. *The Journal of pediatrics* **149**, 809–816.
- 420 OPSOMER, J.D., CLAESKENS, G., RANALLI, M.G., KAUEMANN, G. & BREIDT, F. (2008). Non-
421 parametric small area estimation using penalized spline regression. *Journal of the Royal Statistical*
422 *Society: Series B (Statistical Methodology)* **70**, 265–286.

- 423 PFEFFERMANN, D. et al. (2013). New important developments in small area estimation. *Statistical Science*
424 **28**, 40–68.
- 425 PRATESI, M., RANALLI, M.G. & SALVATI, N. (2008). Semiparametric m-quantile regression for estimating
426 the proportion of acidic lakes in 8-digit hucs of the northeastern us. *Environmetrics* **19**, 687–701.
- 427 PRATESI, M., RANALLI, M.G. & SALVATI, N. (2009). Nonparametric m-quantile regression using penalised
428 splines. *Journal of Nonparametric Statistics* **21**, 287–304.
- 429 RAO, J. & MOLINA, I. (2015). *Small Area Estimation*. New Jersey: John Wiley & Sons, Inc, 2nd edn.
- 430 RAO, J.N., SINHA, S.K. & DUMITRESCU, L. (2014). Robust small area estimation under semi-parametric
431 mixed models. *Canadian Journal of Statistics* **42**, 126–141.
- 432 RUPPERT, D. (2002). Selecting the number of knots for penalized splines. *Journal of computational and*
433 *graphical statistics* **11**, 735–757.
- 434 RUPPERT, D., WAND, M.P. & CARROLL, R.J. (2003). *Semiparametric regression*. 12, Cambridge university
435 press.
- 436 SALVATI, N., RANALLI, M. & PRATESI, M. (2011). Small area estimation of the mean using non-parametric
437 m-quantile regression: a comparison when a linear mixed model does not hold. *Journal of Statistical*
438 *Computation and Simulation* **81**, 945–964.
- 439 SHOKOOHI, F. & TORABI, M. (2018). Semi-parametric small-area estimation by combining time-series and
440 cross-sectional data methods. *Australian & New Zealand Journal of Statistics* **60**, 323–342.
- 441 TORABI, M. (2011). Small area estimation using survey weights with functional measurement error in the
442 covariate. *Australian & New Zealand Journal of Statistics* **53**, 141–155.
- 443 TORABI, M. & SHOKOOHI, F. (2015). Non-parametric generalized linear mixed models in small area
444 estimation. *Canadian Journal of Statistics* **43**, 82–96.
- 445 TORKASHVAND, E., JAFARI JOZANI, M. & TORABI, M. (2015). Pseudo-empirical bayes estimation of small
446 area means based on james–stein estimation in linear regression models with functional measurement
447 error. *Canadian Journal of Statistics* **43**, 265–287.
- 448 YBARRA, L.M. & LOHR, S.L. (2008). Small area estimation when auxiliary information is measured with
449 error. *Biometrika* **95**, 919–931.

Table 1. Estimates (and standard errors) of the model parameters for the proposed model in the case of three forms for spline (linear, quadratic, exponential 1). True values in the case of linear model: $b_0 = b_1 = \sigma_\nu^2 = \sigma_e^2 = 1, \sigma_\eta^2 = 2$; in the case of quadratic model: $b_0 = b_1 = 0.4, b_2 = -0.65, \sigma_\nu^2 = 0.1, \sigma_e^2 = 0.3, \sigma_\eta^2 = 0.6$; in the case of exponential 1 model: $b_0 = b_1 = 1, b_2 = -0.7, \sigma_\nu^2 = 0.1, \sigma_e^2 = 0.3, \sigma_\eta^2 = 0.6$.

True model	b_0	b_1	σ_γ^2	σ_ν^2	σ_e^2	σ_η^2
Linear	0.94(0.89)	1.00(0.04)	$10^{-4}(5 \times 10^{-4})$	0.55(0.30)	1.00(0.13)	1.99(0.26)
Quadratic	-1.07(0.07)	0.13(0.06)	0.10(0.01)	$0.10(7 \times 10^{-13})$	0.31(0.04)	0.59(0.07)
Exponential 1	-0.97(0.08)	-0.73(0.07)	0.18(0.02)	$0.10(5 \times 10^{-13})$	0.30(0.04)	0.60(0.08)

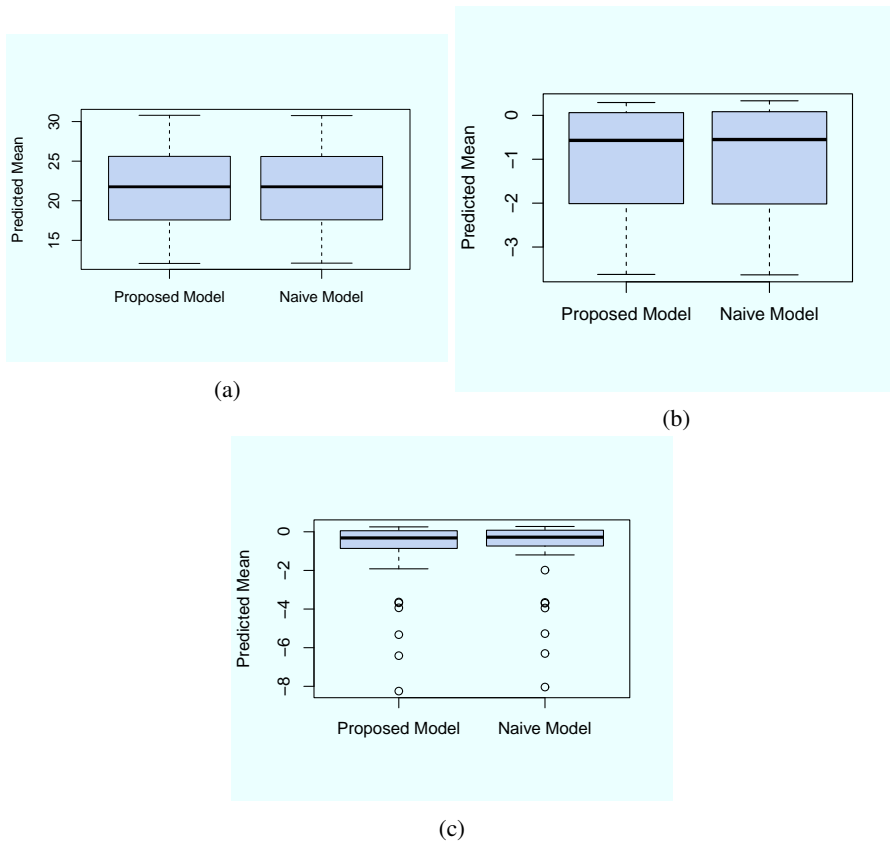


Figure 1. Boxplots of PEB predictions of small area means for the proposed and naive models in the case of (a) linear, (b) quadratic, (c) exponential 1.

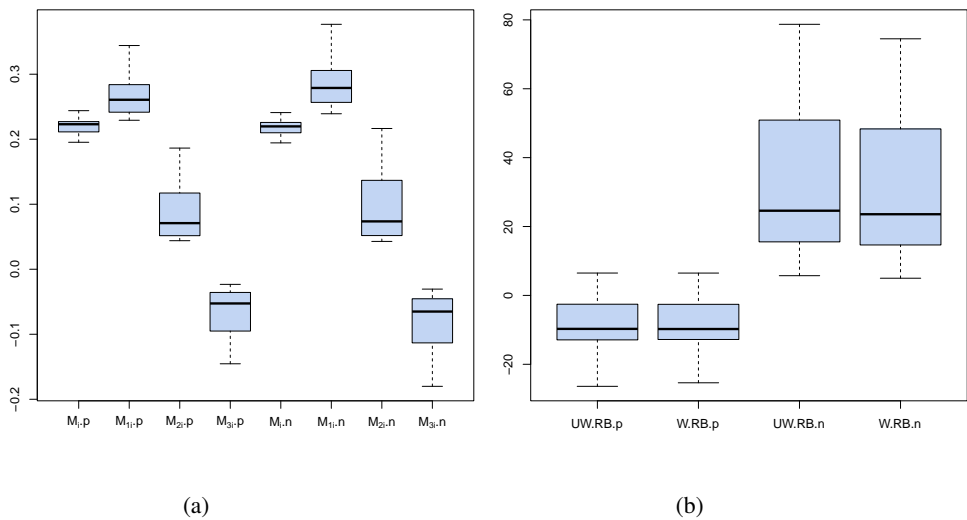


Figure 2. (a) Boxplots of $EMSPE$ of PEB predictors and its components for the proposed model ($M_{i.p}, M_{1i.p}, M_{2i.p}, M_{3i.p}$) and naive model ($M_{i.n}, M_{1i.n}, M_{2i.n}, M_{3i.n}$) in the case of linear form for spline; (b) boxplots of percent relative bias of jackknife estimators of unweighted (UW.RB.p) and weighted (W.RB.p) $MSPE$ for the proposed model, unweighted (UW.RB.n) and weighted (W.RB.n) jackknife $MSPE$ estimation for the naive model in the case of linear form for spline.

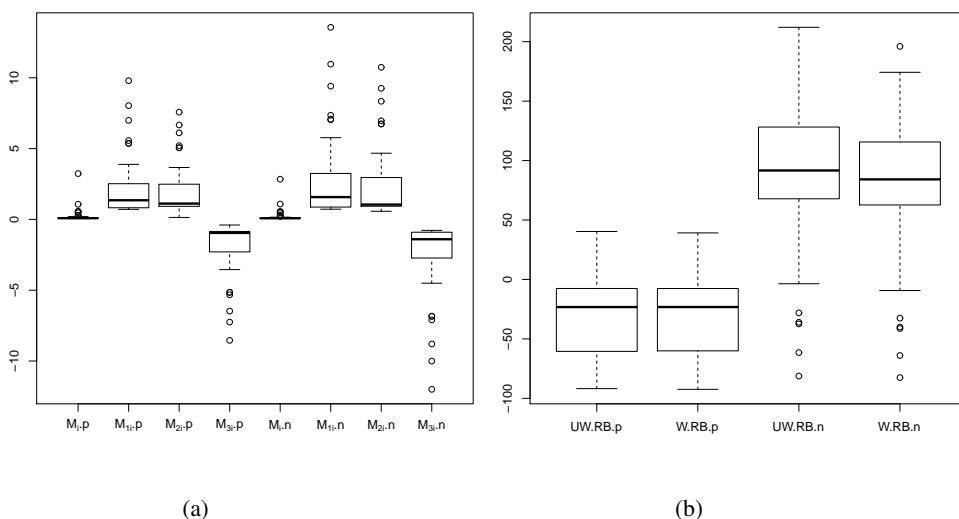


Figure 3. (a) Boxplots of $EMSPE$ of PEB predictors and its components for the proposed model ($M_{i.p}, M_{1i.p}, M_{2i.p}, M_{3i.p}$) and naive model ($M_{i.n}, M_{1i.n}, M_{2i.n}, M_{3i.n}$) in the case of quadratic form for spline; (b) boxplots of percent relative bias of jackknife estimators of unweighted (UW.RB.p) and weighted (W.RB.p) $MSPE$ for the proposed model, unweighted (UW.RB.n) and weighted (W.RB.n) jackknife $MSPE$ estimation for the naive model in the case of quadratic form for spline.

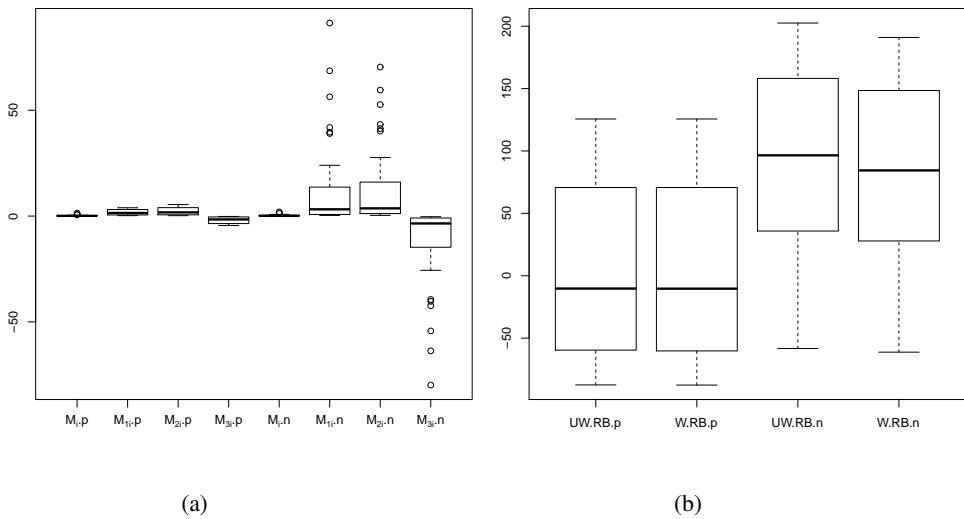


Figure 4. (a) Boxplots of $EMSPE$ of PEB predictors and its components for the proposed model ($M_{i.p}, M_{1i.p}, M_{2i.p}, M_{3i.p}$) and naive model ($M_{i.n}, M_{1i.n}, M_{2i.n}, M_{3i.n}$) in the case of exponential 1 form for spline; (b) boxplots of percent relative bias of jackknife estimators of unweighted (UW.RB.p) and weighted (W.RB.p) $MSPE$ for the proposed model, unweighted (UW.RB.n) and weighted (W.RB.n) jackknife $MSPE$ estimation for the naive model in the case of exponential 1 form for spline.

Table 2. Estimates (and standard errors) of the model parameters for the proposed model in the case of three forms for spline (linear, quadratic, exponential 1). True values in the case of linear model: $b_0 = b_1 = \sigma_\nu^2 = \sigma_e^2 = 1$; in the case of quadratic model: $b_0 = b_1 = 0.4, b_2 = -0.65, \sigma_\nu^2 = 0.1, \sigma_e^2 = 0.3$; in the case of exponential 1 model: $b_0 = b_1 = 1, b_2 = -0.7, \sigma_\nu^2 = 0.1, \sigma_e^2 = 0.3$.

True model	True value of σ_η^2	b_0	b_1	σ_γ^2	σ_ν^2	σ_e^2	σ_η^2
Linear	1	1.00 (0.20)	1.01 (0.14)	$2 \times 10^{-3} (6 \times 10^{-3})$	0.68 (0.26)	1.00 (0.13)	1.01 (0.12)
	3.5	1.00 (0.24)	1.05 (0.22)	$10^{-3} (4 \times 10^{-3})$	0.37 (0.36)	1.00 (0.13)	3.52 (0.43)
	5	0.99 (0.28)	1.08 (0.30)	$10^{-3} (4 \times 10^{-3})$	0.29 (0.36)	1.00 (0.13)	5.03 (0.62)
Quadratic	0.2	-1.07 (0.06)	0.12 (0.04)	0.10 (0.01)	$0.10 (3 \times 10^{-3})$	0.30 (0.03)	0.19 (0.02)
	1	-1.07 (0.06)	0.12 (0.06)	0.12 (0.02)	$0.10 (7 \times 10^{-13})$	0.30 (0.04)	0.93 (0.12)
	1.5	-1.06 (0.06)	0.11 (0.07)	0.14 (0.03)	$0.10 (6 \times 10^{-13})$	0.30 (0.03)	1.45 (0.16)
Exponential 1	0.2	-0.97 (0.07)	-0.73 (0.05)	0.17 (0.01)	$0.10 (5 \times 10^{-13})$	0.30 (0.03)	0.19 (0.02)
	1	-0.97 (0.09)	-0.74 (0.09)	0.19 (0.03)	$0.10 (6 \times 10^{-13})$	0.30 (0.03)	1.00 (0.13)
	1.5	-0.97 (0.09)	-0.75 (0.10)	0.18 (0.03)	$0.10 (6 \times 10^{-13})$	0.30 (0.03)	1.48 (0.18)

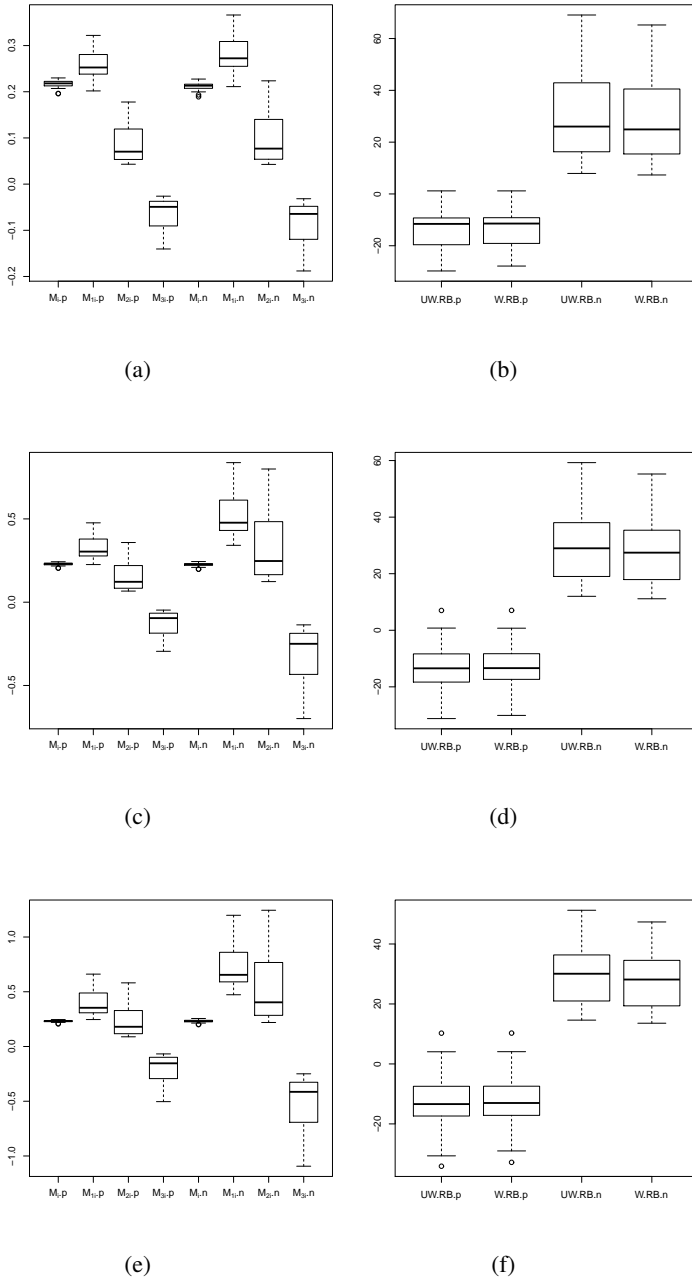


Figure 5. Boxplots of $EMSPE$ of PEB predictors and its components for the proposed model ($M_{i.p}, M_{1i.p}, M_{2i.p}, M_{3i.p}$) and naive model ($M_{i.n}, M_{1i.n}, M_{2i.n}, M_{3i.n}$) in the case of linear form for spline with (a) $\sigma_\eta^2 = 1.00$, (c) $\sigma_\eta^2 = 3.50$, (e) $\sigma_\eta^2 = 5.00$; boxplots of percent relative bias of jackknife estimators of unweighted ($UW.RB.p$) and weighted ($W.RB.p$) $MSPE$ for the proposed model, unweighted ($UW.RB.n$) and weighted ($W.RB.n$) jackknife $MSPE$ estimation for the naive model in the case of linear form for spline with (b) $\sigma_\eta^2 = 1.00$, (d) $\sigma_\eta^2 = 3.50$, (f) $\sigma_\eta^2 = 5.00$.

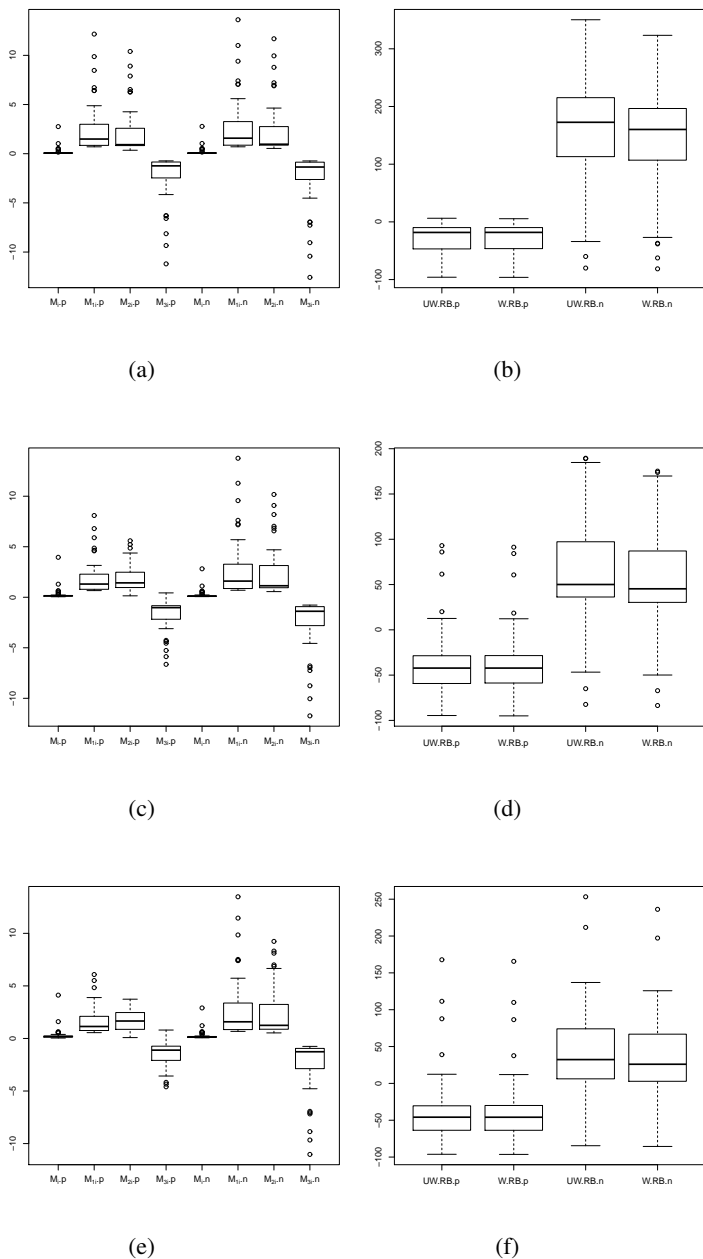


Figure 6. Boxplots of $EMSPE$ of PEB predictors and its components for the proposed model ($M_{i.p}, M_{1i.p}, M_{2i.p}, M_{3i.p}$) and naive model ($M_{i.n}, M_{1i.n}, M_{2i.n}, M_{3i.n}$) in the case of quadratic form for spline with (a) $\sigma_\eta^2 = 0.20$, (c) $\sigma_\eta^2 = 1.00$, (e) $\sigma_\eta^2 = 1.5$; boxplots of percent relative bias of jackknife estimators of unweighted ($UW.RB.p$) and weighted ($W.RB.p$) $MSPE$ for the proposed model, unweighted ($UW.RB.n$) and weighted ($W.RB.n$) jackknife $MSPE$ estimation for the naive model in the case of quadratic form for spline with (b) $\sigma_\eta^2 = 0.20$, (d) $\sigma_\eta^2 = 1.00$, (f) $\sigma_\eta^2 = 1.50$.

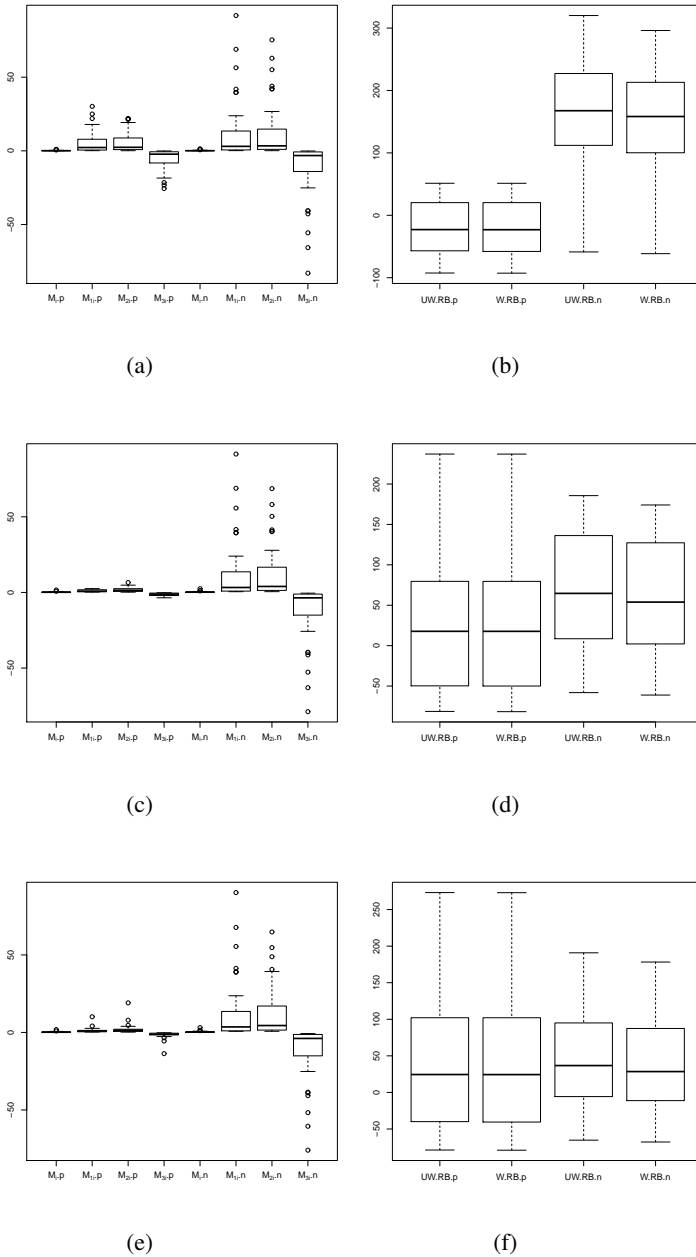


Figure 7. Boxplots of $EMSPE$ of PEB predictors and its components for the proposed model $(M_{i.p}, M_{1i.p}, M_{2i.p}, M_{3i.p})$ and naive model $(M_{i.n}, M_{1i.n}, M_{2i.n}, M_{3i.n})$ in the case of exponential 1 form for spline with (a) $\sigma_\eta^2 = 0.20$, (c) $\sigma_\eta^2 = 1.00$, (e) $\sigma_\eta^2 = 1.50$; boxplots of percent relative bias of jackknife estimators of unweighted (UW.RB.p) and weighted (W.RB.p) $MSPE$ for the proposed model, unweighted (UW.RB.n) and weighted (W.RB.n) jackknife $MSPE$ estimation for the naive model in the case of exponential 1 form for spline with (b) $\sigma_\eta^2 = 0.20$, (d) $\sigma_\eta^2 = 1.00$, (f) $\sigma_\eta^2 = 1.50$.

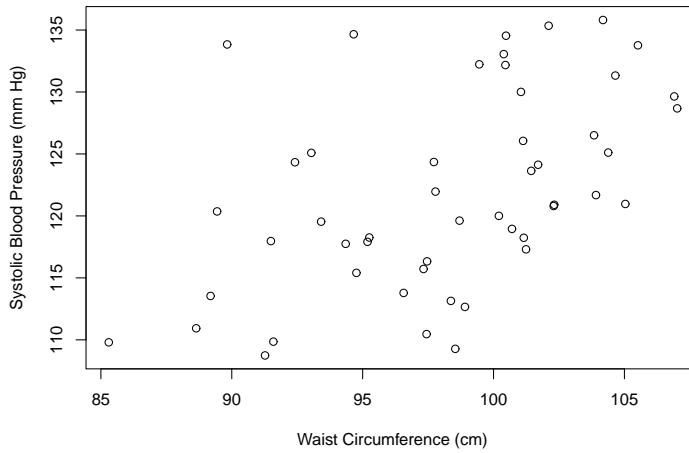


Figure 8. Average systolic blood pressure versus average waist circumference for some predefined groups (sex-age-race and ethnicity) based on US NHANES 2013– 2014 data.

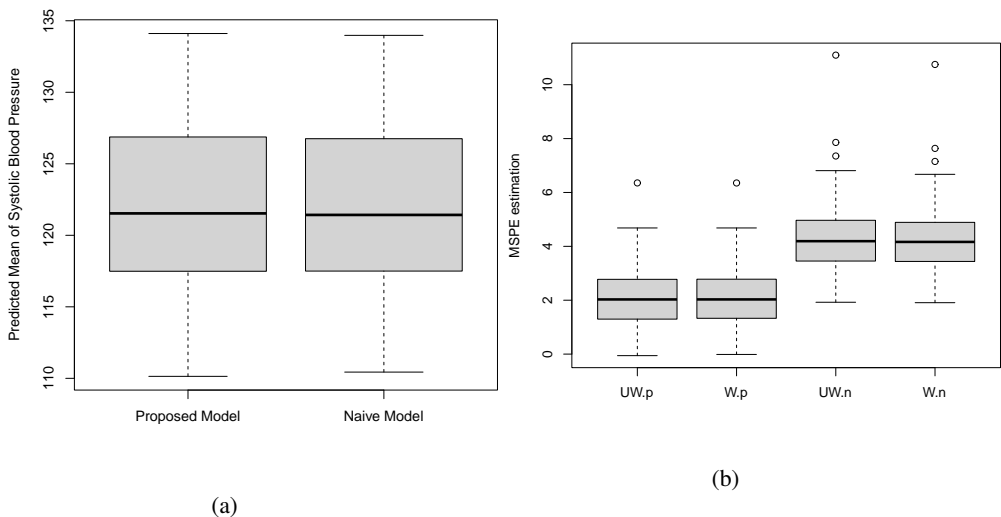


Figure 9. (a) Boxplots of PEB predictors of small area blood pressure means for the proposed and naive models; (b) boxplots of unweighted (UW.p) and weighted (W.p) jackknife estimates of $MSPE$ of small area blood pressure mean predictors for the proposed model; boxplots of unweighted (UW.n) and weighted (W.n) jackknife estimates of $MSPE$ of small area blood pressure mean predictors for the naive model. The boxplots in parts (a) and (b) are for some predefined groups (sex-age-race and ethnicity) based on US NHANES 2013– 2014 data.

Appendix A Supplementary Table

Table A1. Variable definitions

Variable	Definition
m	Number of areas
N_i	Population size in i -th small area
n_i	Sample size in i -th small area
n_T	Total sample size
y_{ij}	Response variable for j -th unit at the i -th small area
\mathbf{y}_i	Vector of response variable in i -th area
\mathbf{y}	Vector of response variable
x_i	True value of the covariate in i -th area
\mathbf{x}_i	Vector of true value of the covariates of fixed part of the model
\mathbf{X}	Matrix of true value of the covariates of fixed part of the model
\mathbf{z}_i	Vector of true value of the covariates of P-spline part of the model
\mathbf{Z}	Matrix of true value of the covariates of P-spline part of the model
w_{ij}	Observed value of the covariate for j -th unit at the i -th small area
$\boldsymbol{\nu} = (\nu_1, \dots, \nu_m)$	Vector of area-level random effects
\mathbf{d}_{ij}	Vector that shows (i, j) -th sample belongs to which area
\mathbf{D}	Matrix that shows each sample belongs to which area
e_{ij}	Random error for j -th unit at the i -th small area
\mathbf{e}_i	Vector of random errors in i -th area
\mathbf{e}	Matrix of random errors
η_{ij}	Measurement error for j -th unit at the i -th small area
$\mathbf{b} = (b_0, \dots, b_p)$	Vector of regression coefficients of fixed part of the model
$\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_k)$	Vector of regression coefficients of P-spline part of the model
p	Degree of spline
(τ_1, \dots, τ_k)	Set of knots
σ_{ν}^2	Variance of the area-level random effects
σ_{γ}^2	Variance of the random effects of P-spline part of the model
σ_{ϵ}^2	Variance of the random error
σ_{η}^2	Variance of the measurement error
\mathbf{V}_{η}	Variance-covariance matrix of \mathbf{Y}
$\boldsymbol{\Sigma}_{\nu}$	Variance-covariance matrix of $\boldsymbol{\nu}$
$\boldsymbol{\Sigma}_{\gamma}$	Variance-covariance matrix of $\boldsymbol{\gamma}$
$\boldsymbol{\Sigma}_{\mathbf{e}}$	Variance-covariance matrix of \mathbf{e}
$\boldsymbol{\Sigma}_{\eta}$	Variance-covariance matrix of $\boldsymbol{\eta}$
θ_i	Mean response variable of i -th small area
θ_i^B	Best predictor of i -th small area
θ_i^{PB}	Pseudo-best predictor of i -th area
θ_i^{PEB}	Pseudo-empirical best predictor of i -th small area
\tilde{x}_i	Estimate of x_i when all the parameters are known
\hat{x}_i	Estimate of x_i when all the parameters are estimated
$\mathbf{h}_i = (1, \bar{W}_i, \dots, \bar{W}_i^p)^{\top}$	Vector of the observed value of the covariates mean
R	Number of simulation runs in simulation study
$m\text{spe}_J$	Unweighted jackknife mean squared prediction error estimator
$m\text{spe}_{JW}$	Weighted jackknife mean squared prediction error estimator
RB	Relative Bias