



Research article

Bayesian analysis of random effects panel interval-valued data models

Dengke Xu*, Linlin Shen, Yuanyang Tangzhu and Shiqi Ke

School of Economics, Hangzhou Dianzi University, Hangzhou 310018, China

* **Correspondence:** Email: xudengke1983@163.com.

Abstract: In the era of big data, interval-valued data is quite common in real life and can be used to describe the uncertainty of variables. In this paper, we introduced random effects panel interval-valued data models based on the center and range method and constructed a Bayesian method for the models, including estimation and prediction. Some simulation studies indicate that the proposed Bayesian method performs well. Finally, our proposed panel interval-valued data Bayesian models were applied in forecasting of the Air Quality Index, and the experimental evaluation of actual data sets shows the advantages and the performance of our proposed models.

Keywords: panel data; interval-valued data; Bayesian estimation; center and range method

1. Introduction

Nowadays, big data, artificial intelligence, 5G, and other technologies have gradually entered people's vision, and have a huge impact on our lives. These also promote the total amount of data to show exponential growth. In the face of massive data, how to quickly and accurately obtain, process, mine, and integrate the required information becomes particularly important. According to different actual situations and purpose requirements, the data type is no longer limited to the definite point value data, but also gradually evolves to the uncertain data. As time goes on, many scholars at home and abroad have proposed interval numbers to describe uncertain data. Generally speaking, interval-valued data [1] arise due to one of the following two reasons: (i) imprecise observations of quantities, resulting in the translation of the measured value into an interval of possible values, and (ii) information aggregation. The past three decades have witnessed enormous developments in statistical methods for interval-valued data. In particular, methods of regression focusing on interval-valued data have been extensively developed. For instance, Billard and Diday [2] proposed the center method (CM) for linear regression analysis of the interval-valued data, assuming that the interval-valued data had a uniform distribution. Billard and Diday [3] proposed the min-max method, taking the lower and upper limits of interval data as special point values, and establishing linear regression models for each of them. Billard and

Diday [4] proposed the binary center and range method (BCRM), in which the explanatory variable contains information about both the center and range in the regression process. Then, LimaNeto and De Carvalho [5] proposed the centre and range method (CRM), which implied that the intervals were converted to centers and ranges. Based on the CRM, Lima Neto and Carvalho [6] introduced the constrained center and range method (CCRM), which added non-negative constraints of range coefficients. Giordani [7] adapted the CRM and presented a Lasso-based method for the interval-valued regression model. In addition, some scholars tried to establish interval-valued data regression models by other methods instead of the CRM. Maia and De Carvalho [8] proposed a multi-layer perceptron neural network based on interval-valued data and the Holt exponential smoothing method. Souza et al. [9] considered a parametrized approach that automatically extracted the best reference points from intervals. Kong et al. [10] proposed an interval local linear method (ILLM) to fit a regression model with interval-valued explanatory and response variables, which has no restrictions on the form of the regression function. Kong and Gao [11] studied the method of moments (MM) estimation method for interval-valued regression models. However, there is almost no literature studying the Bayesian modeling problem of interval panel data models.

As is known to all, a panel data set offers a certain number of advantages over traditional pure cross-section or pure time series data sets and statistical models combining cross-section and time series real-valued data become increasingly popular in economic research. For example, Nuroglu and Kunst [12] discussed the impact of exchange rate fluctuations on international trade flows by using panel data analysis and fuzzy data analysis methods. He et al. [13] constructed an interval slacks-based measure (SBM) of non-expected output, which analyzed China's environmental technology efficiency based on panel data of various provinces. But so far, the data in panel data models are all real-valued [14–19], so it is very necessary to build a panel data model for interval-valued data. Recently, Ji et al. [20] introduced a panel data regression model for interval-valued data and constructed three kinds of panel interval-valued data regression models, and this is the first attempt to discuss panel interval-valued data models.

In addition, due to the recent dramatic evolution in advanced computational technologies, Bayesian inference has also received a lot of attention in recent years. In Bayesian regression fields, Park and Casella [21] proposed a Bayesian Lasso method for linear models. By using spline approximation, Xu and Zhang [22] introduced a Bayesian method for the partially linear model with heteroscedasticity based on the variance modeling technique. Castillo et al. [23] studied a high-dimensional linear regression with a sparse prior, which is a mixture of point masses at zero and continuous distributions. Pfarrhofer and Piribauer [24] proposed two shrink age priors to make Bayesian variable selection for high-dimensional spatial autoregressive models. Wang and Tang [25] considered Bayesian inference on a quantile regression model in the presence of nonignorable missing covariates. Zhang et al. [26] considered Bayesian quantile regression analysis for semiparametric mixed-effects double regression models based on the asymmetric Laplace distribution for the errors. Tang et al. [27] for the first time used Elastic Net, a penalized method, for Bayesian quantile regression of panel data and derived the posterior distributions of all parameters based on the asymmetric Laplace prior distributions, and then constructed Gibbs sampling. Tao et al. [28] proposed a Bayesian adaptive Lasso quantile regression method based on asymmetric exponential power distribution and applied the method to panel data. However, few works are devoted to constructing the Bayesian framework for interval-valued data. Zhang et al. [29] proposed the Bayesian nonparametric regression models by assuming that the upper and lower

of the interval were distributed as an asymmetric Laplace distribution. Xu and Qin [30] extended the CRM for interval-valued regression models to the Bayesian framework for the first time, and proposed a bivariate Bayesian regression model based on the CRM with known and unknown covariance matrices.

However, to the best of our knowledge, there is little work done for constructing the Bayesian framework for panel interval-valued data. Hence, a Bayesian model for random effects panel interval-valued data is developed on the basis of the center and range method in this paper, and is compared with the Bayesian estimation based on the center method and the Bayesian estimation based on the minimum and maximum method. Finally, our proposed panel interval-valued data models are applied in forecasting of the Air Quality Index.

The outline of the paper is as follows. In Section 2, we introduce a Bayesian model of random effects panel interval-valued data based on the center and range method, and give the likelihood based on this model. In Section 3, the prior distribution and Bayesian posterior inference of the model are given in detail. In Section 4, the results of parameter estimation are obtained through simulation and compared with other methods to illustrate the feasibility of the proposed method. In Section 5, we apply the model to real data, and in Section 6, the paper is concluded with a brief discussion.

2. Random effects panel interval-valued data models

For panel interval-valued data set $S = \{(X_{it}, y_{it}) | i = 1, 2, \dots, n; t = 1, 2, \dots, T\}$, $y_{it} = [y_{it}^l, y_{it}^u]$ is assumed as the observed interval-valued dependent variables, where the superscripts l and u represent the lower and upper bounds of the interval, respectively, and let $X_{it} = (x_{it1}, x_{it2}, \dots, x_{itp})^T$ be $p \times 1$ interval-valued independent vectors with $x_{itj} = [x_{itj}^l, x_{itj}^u]$, $i = 1, 2, \dots, n; t = 1, 2, \dots, T; j = 1, 2, \dots, p$.

In the following, we consider Bayesian analysis of random effects panel interval-valued data models:

$$\begin{aligned} y_{it}^c &= \alpha_i^c + X_{it}^c \beta^c + \varepsilon_{it}^c, \\ y_{it}^r &= \alpha_i^r + X_{it}^r \beta^r + \varepsilon_{it}^r, \end{aligned} \quad (2.1)$$

where $y_{it}^c = \frac{y_{it}^u + y_{it}^l}{2}$, $y_{it}^r = \frac{y_{it}^u - y_{it}^l}{2}$, $X_{it}^c = \frac{X_{it}^u + X_{it}^l}{2}$, $X_{it}^r = \frac{X_{it}^u - X_{it}^l}{2}$. $\beta^c = (\beta_1^c, \dots, \beta_p^c)^T$ and $\beta^r = (\beta_1^r, \dots, \beta_p^r)^T$ are vectors of p -dimensional unknown parameters. In addition, ε_{it}^c and ε_{it}^r are mutually independent and identically distributed normal random variables with zero mean and variances σ_c^2 and σ_r^2 , respectively. α_i^c and α_i^r denote the random effects associated with individual i and assume that the random effects $\alpha_i^c \sim N(0, \phi_c^2)$, $\alpha_i^r \sim N(0, \phi_r^2)$. Then for time t , model (2.1) can be written as

$$\begin{aligned} Y_t^c &= \alpha^c + X_t^c \beta^c + \varepsilon_t^c, \\ Y_t^r &= \alpha^r + X_t^r \beta^r + \varepsilon_t^r. \end{aligned} \quad (2.2)$$

Then for $s \in \{c, r\}$, the center and range models take on the following form:

$$Y_t^s = \alpha^s + X_t^s \beta^s + \varepsilon_t^s, \quad (2.3)$$

where $Y_t^s = (y_{1t}^s, y_{2t}^s, \dots, y_{nt}^s)^T$, $\alpha^s = (\alpha_1^s, \alpha_2^s, \dots, \alpha_n^s)^T$, $X_t^s = (X_{1t}^s, X_{2t}^s, \dots, X_{nt}^s)^T$, $\beta^s = (\beta_1^s, \beta_2^s, \dots, \beta_p^s)^T$, $\varepsilon_t^s = (\varepsilon_{1t}^s, \varepsilon_{2t}^s, \dots, \varepsilon_{nt}^s)^T$.

For convenience, matrix notation and vector notation are used to represent variables and models. Let $Y^s = ((Y_1^s)^T, (Y_2^s)^T, \dots, (Y_T^s)^T)^T$, $X^s = ((X_1^s)^T, (X_2^s)^T, \dots, (X_T^s)^T)^T$, $\varepsilon^s = ((\varepsilon_1^s)^T, (\varepsilon_2^s)^T, \dots, (\varepsilon_T^s)^T)^T$,

$\mathbf{D}^s = \mathbf{1}_T^s \otimes \mathbf{I}_n^s$, where “ \otimes ” denotes the Kronecker product, $\mathbf{1}_T^s$ is a $T \times 1$ dimensional vector with all elements 1, and \mathbf{I}_n^s is an $n \times n$ dimensional unit matrix. Then, model (2.3) can also be written in matrix form,

$$\mathbf{Y}^s = \mathbf{D}^s \boldsymbol{\alpha}^s + \mathbf{X}^s \boldsymbol{\beta}^s + \boldsymbol{\varepsilon}^s, \quad (2.4)$$

where $\boldsymbol{\alpha}^s \sim N(0, \phi_s^2 \mathbf{I}_n)$, $\boldsymbol{\varepsilon}^s \sim N(0, \sigma_s^2 \mathbf{I}_N)$, $N = n \times T$.

Based on the above model, the likelihood function of the model parameters is defined as follows:

$$L(\mathbf{Y}^s, \boldsymbol{\alpha}^s | \mathbf{X}^s, \boldsymbol{\beta}^s, \sigma_s^2, \phi_s^2) = (2\pi)^{-\frac{N+n}{2}} (\sigma_s^2)^{-\frac{N}{2}} (\phi_s^2)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma_s^2} (\mathbf{Y}^s - \mathbf{D}^s \boldsymbol{\alpha}^s - \mathbf{X}^s \boldsymbol{\beta}^s)^T (\mathbf{Y}^s - \mathbf{D}^s \boldsymbol{\alpha}^s - \mathbf{X}^s \boldsymbol{\beta}^s) - \frac{(\boldsymbol{\alpha}^s)^T \boldsymbol{\alpha}^s}{2\phi_s^2} \right\}. \quad (2.5)$$

3. Bayesian analysis of the models

To estimate the unknown parameters $\boldsymbol{\beta}^s$, σ_s^2 , and ϕ_s^2 , we implement a Bayesian approach [31]. Therefore, we should appoint prior distributions for the parameters of the models. First, we assume $\boldsymbol{\beta}^s$ has normal prior distributions, i.e., $\boldsymbol{\beta}^s \sim N(\boldsymbol{\beta}_0^s, \boldsymbol{\Sigma}_0^s)$. It is assumed that the hyperparameters $\boldsymbol{\beta}_0^s$, $\boldsymbol{\Sigma}_0^s$ are known. In addition, the prior information for other unknown parameters is $\sigma_s^2 \sim IG(a_0^s, b_0^s)$, $\phi_s^2 \sim IG(c_0^s, d_0^s)$, where a_0^s , b_0^s , c_0^s , d_0^s are hyperparameters to be given, and “ IG ” means the inverse gamma distribution. Thus, the joint priors of all of the unknown parameter are given by

$$\pi(\boldsymbol{\beta}^s, \sigma_s^2, \phi_s^2) = p(\boldsymbol{\beta}^s) p(\sigma_s^2) p(\phi_s^2). \quad (3.1)$$

From the likelihood function (2.5) and prior distributions (3.1), we can obtain the following theorems and give a brief proof of Theorem 1. The proofs of the other theorems are similar and will not be written in detail here.

Theorem 1. Suppose that the parameter $\boldsymbol{\beta}^s$ follows a normal prior distribution, i.e., $\boldsymbol{\beta}^s \sim N(\boldsymbol{\beta}_0^s, \boldsymbol{\Sigma}_0^s)$, and then the posterior distribution of $\boldsymbol{\beta}^s$ follows the normal distribution,

$$p(\boldsymbol{\beta}^s | \mathbf{Y}^s, \mathbf{X}^s, \boldsymbol{\alpha}^s, \sigma_s^2, \phi_s^2) \sim N(\boldsymbol{\mu}_{\boldsymbol{\beta}^s}, \boldsymbol{\Sigma}_{\boldsymbol{\beta}^s}),$$

where $\boldsymbol{\mu}_{\boldsymbol{\beta}^s} = \boldsymbol{\Sigma}_{\boldsymbol{\beta}^s} \left(\frac{(\mathbf{X}^s)^T (\mathbf{Y}^s - \mathbf{D}^s \boldsymbol{\alpha}^s)}{\sigma_s^2} + (\boldsymbol{\Sigma}_0^s)^{-1} \boldsymbol{\beta}_0^s \right)$, $\boldsymbol{\Sigma}_{\boldsymbol{\beta}^s} = \left(\frac{1}{\sigma_s^2} (\mathbf{X}^s)^T \mathbf{X}^s + (\boldsymbol{\Sigma}_0^s)^{-1} \right)^{-1}$.

The posterior distribution of $\boldsymbol{\beta}^s$ is as follows, and we make a transformation of its form:

$$\begin{aligned} p(\boldsymbol{\beta}^s | \mathbf{Y}^s, \mathbf{X}^s, \boldsymbol{\alpha}^s, \sigma_s^2, \phi_s^2) &\propto L(\mathbf{Y}^s, \boldsymbol{\alpha}^s | \mathbf{X}^s, \boldsymbol{\beta}^s, \sigma_s^2, \phi_s^2) \pi(\boldsymbol{\beta}^s, \sigma_s^2, \phi_s^2) \\ &\propto \exp \left\{ -\frac{1}{2\sigma_s^2} (\mathbf{Y}^s - \mathbf{D}^s \boldsymbol{\alpha}^s - \mathbf{X}^s \boldsymbol{\beta}^s)^T (\mathbf{Y}^s - \mathbf{D}^s \boldsymbol{\alpha}^s - \mathbf{X}^s \boldsymbol{\beta}^s) \right\} \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta}^s - \boldsymbol{\beta}_0^s)^T (\boldsymbol{\Sigma}_0^s)^{-1} \right. \\ &\quad \left. (\boldsymbol{\beta}^s - \boldsymbol{\beta}_0^s) \right\} \\ &\propto \exp \left\{ -\frac{1}{2} (\boldsymbol{\beta}^s)^T \left(\frac{(\mathbf{X}^s)^T \mathbf{X}^s}{\sigma_s^2} + (\boldsymbol{\Sigma}_0^s)^{-1} \right) \boldsymbol{\beta}^s + \left(\frac{(\mathbf{Y}^s - \mathbf{D}^s \boldsymbol{\alpha}^s)^T \mathbf{X}^s}{\sigma_s^2} + (\boldsymbol{\beta}_0^s)^T (\boldsymbol{\Sigma}_0^s)^{-1} \right) \boldsymbol{\beta}^s \right\}. \end{aligned}$$

Obviously, the conditional distribution of $\boldsymbol{\beta}^s$ agrees with the form of the probability density of the multivariate normal distribution, so $\boldsymbol{\beta}^s$ obeys the normal distribution, thus this paper obtains Theorem 1.

Theorem 2. Based on model assumptions, we can obtain the posterior distribution of α^s as follows, that is

$$p(\alpha^s | Y^s, X^s, \beta^s, \sigma_s^2, \phi_s^2) \sim N(\mu_{\alpha^s}, \Sigma_{\alpha^s}),$$

where $\mu_{\alpha^s} = \Sigma_{\alpha^s}^{-1} (\mathbf{D}^s)^T (Y^s - X^s \beta^s)$, $\Sigma_{\alpha^s} = (\frac{(\mathbf{D}^s)^T \mathbf{D}^s}{\sigma_s^2} + \frac{\mathbf{I}_{n^s}}{\phi_s^2})^{-1}$.

Theorem 3. Suppose that the prior distribution of σ_s^2 is known, i.e., $\sigma_s^2 \sim IG(a_0^s, b_0^s)$, and then the posterior distribution of σ_s^2 follows the inverse gamma distribution

$$p(\sigma_s^2 | Y^s, X^s, \beta^s, \alpha^s, \phi_s^2) \sim IG(a_*^s, b_*^s),$$

where $a_*^s = \frac{N}{2} + a_0^s$, $b_*^s = \frac{1}{2} (Y^s - \mathbf{D}^s \alpha^s - X^s \beta^s)^T (Y^s - \mathbf{D}^s \alpha^s - X^s \beta^s) + b_0^s$.

Theorem 4. Suppose that the prior distribution of ϕ_s^2 is known, i.e., $\phi_s^2 \sim IG(c_0^s, d_0^s)$, and then the posterior distribution of ϕ_s^2 follows the inverse gamma distribution

$$p(\phi_s^2 | Y^s, X^s, \beta^s, \alpha^s, \sigma_s^2) \sim IG(c_*^s, d_*^s),$$

where $c_*^s = \frac{n}{2} + c_0^s$, $d_*^s = \frac{(\alpha^s)^T \alpha^s}{2} + d_0^s$.

From Theorems 1 to 4, the posteriors of β^s , σ_s^2 , ϕ_s^2 are all familiar distributions and can be sampled directly. The specific algorithm is shown in Table 1. Thus based on the above MCMC algorithm, a converged posterior sample can be collected and this sample is noted as $\theta^{(sim)} = (\beta^{s(sim)}, \sigma_s^{2(sim)}, \phi_s^{2(sim)})$, $sim = 1, 2, \dots, M$, $M < Sim$. As such, the posterior estimation of the parameters $(\hat{\beta}^s, \hat{\sigma}_s^2, \hat{\phi}_s^2)$ can be respectively estimated as follows:

$$\hat{\beta}^s = \frac{1}{M} \sum_{sim=1}^M \beta^{s(sim)}, \hat{\sigma}_s^2 = \frac{1}{M} \sum_{sim=1}^M \sigma_s^{2(sim)}, \hat{\phi}_s^2 = \frac{1}{M} \sum_{sim=1}^M \phi_s^{2(sim)}.$$

During the implementation process, the algorithm is repeated 100 times, and the final parameter estimation is the result based on the 100 times average.

Table 1. The algorithm of Bayesian estimation for the unknown parameters.

Algorithm: The specific algorithm is given for unknown parameters $\theta = (\beta^s, \sigma_s^2, \phi_s^2)$
Input: The initial value $\theta^{(0)} = (\beta^{s(0)}, \sigma_s^{2(0)}, \phi_s^{2(0)})$ is given, and the number of iterations of the sampling algorithm is Sim .
Output: Posterior sample sequence $(\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(Sim)})$.
for $sim \leftarrow 1$ to Sim do:
1. Sampling $\alpha^s Y^s, X^s, \beta^s, \sigma_s^2, \phi_s^2 \sim N(\mu_{\alpha^s}, \Sigma_{\alpha^s})$;
2. Sampling $\beta^s Y^s, X^s, \alpha^s, \sigma_s^2, \phi_s^2 \sim N(\mu_{\beta^s}, \Sigma_{\beta^s})$;
3. Sampling $\sigma_s^2 Y^s, X^s, \beta^s, \alpha^s, \phi_s^2 \sim IG(a_*^s, b_*^s)$;
4. Sampling $\phi_s^2 Y^s, X^s, \beta^s, \alpha^s, \sigma_s^2 \sim IG(c_*^s, d_*^s)$;
End

4. Simulation study

In this section, we investigate the performance of the proposed model and Bayesian estimation method via Monte Carlo simulation. We compare Bayesian analysis of random effects panel interval-valued data models based on the center and range method (BCRM) with Bayesian analysis of random

effects panel interval-valued data models based on the center method (BCM) and Bayesian analysis of random effects panel interval-valued data based on the minimum and maximum method (BMMM). In order to demonstrate the quality of Bayesian estimation and prediction, Section 4.1 introduces the measurement methods. Section 4.2 shows the data generation process and parameter details. Then we present all simulation results and conclusions in Section 4.3.

4.1. Measurements

There are three measurements for evaluating the performances of different models:

1) The upper and lower bounds root mean-square errors, i.e., $RMS E_U$ and $RMS E_L$ [5] are

$$RMS E_U = \sqrt{\frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^T (\hat{y}_{it}^u - y_{it}^u)^2},$$

$$RMS E_L = \sqrt{\frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^T (\hat{y}_{it}^l - y_{it}^l)^2}.$$

2) The root mean-square error $RMS E_H$ [32] is

$$RMS E_H = \sqrt{\frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^T (|\hat{y}_{it}^u - y_{it}^u| + |\hat{y}_{it}^l - y_{it}^l|)^2}.$$

3) The rate of different intervals (RI) defined by Hu and He [33] is

$$RI = \frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^T \frac{\omega(y_{it} \cap \hat{y}_{it})}{\omega(y_{it} \cup \hat{y}_{it})},$$

where $\omega(\cdot)$ represents the width of the intervals.

4.2. Data generating process

In this subsection, based on the center and range method (BCRM), four configurations C1, C2, C3, and C4 are generated, as shown in Table 2.

The data sets are generated from the following models:

$$Y^s = D^s \alpha^s + X^s \beta^s + \epsilon^s. \quad (4.1)$$

Referring to Xu and Qin [30], first we fixed the regression coefficients $\beta^c = (8, 7, 6)^T$ and $\beta^r = (4, 6, 5)^T$. Second, X^s is generated from the uniform distribution $U(0, 2)$ and α^s is generated from $N(0, 0.5I_n)$ and fixed for each replication. We generate Y^s from the above model, where $\epsilon^s \sim N(0, \sigma_s^2 I_N)$. In addition, we consider the noninformative prior information type of hyperparameter values for unknown parameters β^s , σ_s^2 , and ϕ_s^2 in the simulation: $\beta_0^s = \mathbf{0}_3$, $\Sigma_0^s = 10 \times I_3$, $a_0^s = b_0^s = c_0^s = d_0^s = 0.01$, where $\mathbf{0}_3$ is a 3-dimensional vector with all elements being 0. Further, we choose n to be 50, 100 and T to be 12, and therefore the sample sizes are $N = n \times T = 600$ and 1200.

Table 2. Four configurations of simulation data.

C1	C2	C3	C4
$X^s \sim U(0, 2)$	$X^s \sim U(0, 2)$	$X^s \sim U(0, 2)$	$X^s \sim U(0, 2)$
$\alpha^s \sim N(0, 0.5I_n)$	$\alpha^s \sim N(0, 0.5I_n)$	$\alpha^s \sim N(0, 0.5I_n)$	$\alpha^s \sim N(0, 0.5I_n)$
$\varepsilon^c \sim N(0, 0.5I_N)$	$\varepsilon^c \sim N(0, 2I_N)$	$\varepsilon^c \sim N(0, 1I_N)$	$\varepsilon^c \sim N(0, 3I_N)$
$\varepsilon^r = \varepsilon^c + e$	$\varepsilon^r = \varepsilon^c + e$	$\varepsilon^r \sim N(0, 1I_N)$	$\varepsilon^r \sim N(0, 3I_N)$
$e \sim N(0, 0.5I_N)$	$e \sim N(0, 0.5I_N)$		

For configurations C1 and C2, we assume that there is a linear relationship hidden in ε^c and ε^r that $\varepsilon^r = \varepsilon^c + e$, where e is the random error generated from $N(0, \sigma^2 I_N)$ and $\sigma^2 = 0.5$. In configurations C3 and C4 it is assumed that there is no linear relationship between ε^c and ε^r .

Based on the above settings and the generated data sets, the preceding proposed algorithm is used to evaluate the Bayesian estimation of unknown parameters based on 100 replications. In order to obtain better and more accurate results, we collect the observation results of the following $J = 2000$ for statistical inference by discarding the first 3000 burn-in iterations.

4.3. Results

Tables 3 and 4 show the Bayesian estimates when $n = 50, 100$ and T is fixed at 12 based on 100 repetitions in all configurations. The accuracy of the estimation is expressed by BIAS and standard deviation SD, and the following conclusions are obtained:

1) Comparing configurations C1, C2, and C4, it can be found that when n is the same, SD of C1 are basically the smallest among the three configurations, while SD of C4 are basically the largest, which indicates that the smaller the variance of the error term, the smaller the estimated error.

2) As the sample size increases, the accuracy of Bayesian estimation for the BCRM gradually gets better. For example, in configuration C1, when $n = 50$ and 100, the SD of β_1^c is 0.0476 and 0.0336, respectively.

In general, the BIAS and SD of all parameters in the four configurations are very small, indicating that the Bayesian estimation effect is good.

Table 3. The results of Bayesian estimation when $n = 50$.

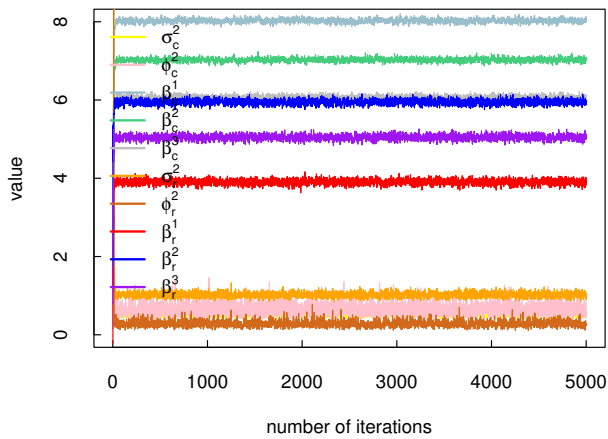
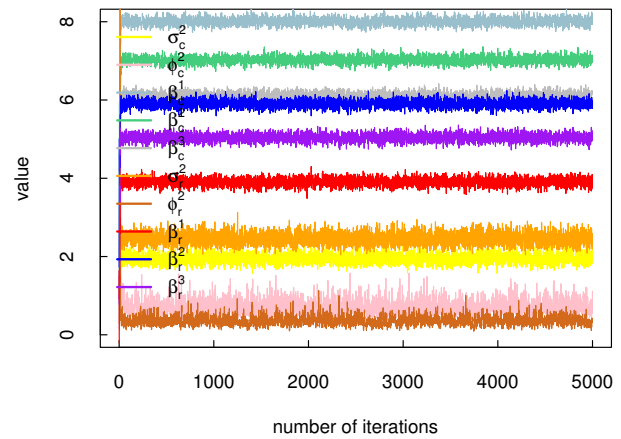
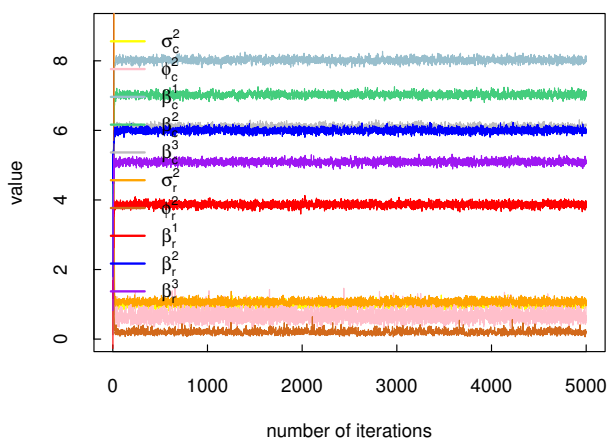
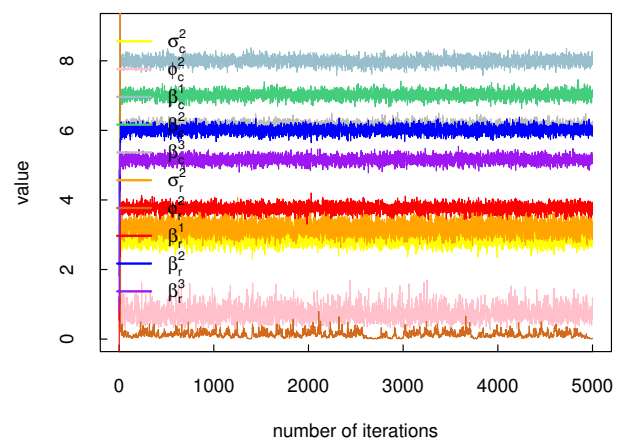
Config.	Para.	β_1^c	β_2^c	β_3^c	σ_c^2	ϕ_c^2	β_1^r	β_2^r	β_3^r	σ_r^2	ϕ_r^2
C1	BIAS	-0.0053	-0.0012	-0.0020	0.0135	0.0098	-0.0110	-0.0080	0.0059	0.0197	0.0151
	SD	0.0476	0.0522	0.0525	0.0319	0.1145	0.0588	0.0602	0.0644	0.0622	0.1228
C2	BIAS	-0.0078	-0.0011	-0.0008	0.0330	-0.0019	-0.0179	-0.0113	0.0089	0.0409	0.0092
	SD	0.0931	0.0964	0.0966	0.1163	0.1350	0.0897	0.0948	0.0932	0.1510	0.1367
C3	BIAS	-0.0063	-0.0010	-0.0016	0.0212	0.0059	-0.0051	-0.0089	0.0028	0.0141	0.0153
	SD	0.0667	0.0707	0.0712	0.0605	0.1211	0.0602	0.0616	0.0706	0.0511	0.1362
C4	BIAS	-0.0090	-0.0014	-0.0005	0.0434	-0.0105	-0.0074	-0.0146	0.0048	0.0235	0.0107
	SD	0.1135	0.1158	0.1160	0.1720	0.1502	0.1021	0.1031	0.1154	0.1476	0.1767

To investigate the convergence of the proposed algorithm, we visualize the values of 5000 iteration updates in a loop for each configuration when $n = 50$, and the results are shown in Figure 1. It can be seen that all the parameters are roughly in a straight line, indicating that the algorithm has a considerable convergence speed.

The next step is to do a predictive study of the models. We compare the Bayesian method based on the center and range method with the other two methods mentioned above, and all results are listed in

Table 4. The results of Bayesian estimation when $n = 100$.

Config.	Para.	β_1^c	β_2^c	β_3^c	σ_c^2	ϕ_c^2	β_1^r	β_2^r	β_3^r	σ_r^2	ϕ_r^2
C1	BIAS	-0.0012	-0.0006	-0.0041	0.0037	0.0044	0.0054	0.0041	-0.0005	0.0032	0.0003
	SD	0.0336	0.0319	0.0366	0.0229	0.0766	0.0476	0.0467	0.0439	0.0381	0.0898
C2	BIAS	-0.0025	-0.0003	-0.0084	0.0040	0.0043	0.0098	-0.0004	-0.0022	0.0004	-0.0010
	SD	0.0631	0.0600	0.0673	0.0879	0.0924	0.0731	0.0682	0.0614	0.0984	0.1042
C3	BIAS	-0.0019	-0.0007	-0.0060	0.0045	0.0043	-0.0003	0.0097	0.0006	0.0057	0.0017
	SD	0.0460	0.0437	0.0494	0.0446	0.0817	0.0465	0.0472	0.0494	0.0393	0.0936
C4	BIAS	-0.0030	0.0001	-0.0102	0.0027	0.0040	-0.0010	0.0142	0.0000	0.0056	0.0013
	SD	0.0761	0.0726	0.0806	0.1311	0.1034	0.0765	0.0788	0.0797	0.1160	0.1197

**(a)** Parameter convergence based on C1**(b)** Parameter convergence based on C2**(c)** Parameter convergence based on C3**(d)** Parameter convergence based on C4**Figure 1.** Parameter convergence with different configurations for $n = 50$.

Tables 5–8, where the standard deviation is in parentheses. 75% of the data is selected as the training set and 25% of the data as the test set. Each case is repeated 100 times and the results of the prediction error are averaged. The following conclusions can be drawn:

1) When n is the same, the BCRM is better than the BCM and the BMMM. For example, in C1 configuration, when $n = 50$, the $RMS E_L$, $RMS E_U$, and $RMS E_H$ of the BCRM are 0.723, 1.585, and 1.631, respectively, and RI is 0.931. The $RMS E_L$, $RMS E_U$, $RMS E_H$, and RI of the BCM and BMMM methods are 6.531, 6.664, 7.094, 0.714 and 6.296, 2.360, 6.320, 0.768, respectively. The first three data are larger than in the BCRM, while RI is smaller than in the BCRM. This shows that the BCRM has better prediction effect than the BCM and the BMMM.

2) Also in the case of $n = 50$, the $RMS E_L$ of the BCRM method in C1 configuration is 0.723 and the standard deviation is 0.041, and that of the BCRM method in C2 configuration is 0.738, 0.044, indicating that the smaller the variance of the random error, the higher the accuracy of the prediction.

3) As the sample size increases, the effectiveness of Bayesian prediction gradually improves. For example, in C2 configuration, the standard deviation of the $RMS E_L$ for the BCRM decreases with increasing the sample size, changing from 0.044 to 0.029.

In general, the prediction effect of the BCRM based on the center and range method is satisfactory in different configurations and different sample sizes.

Table 5. Prediction results of each method in C1 configuration.

n	T	$RMS E_L$			$RMS E_U$		
		BCRM	BCM	BMMM	BCRM	BCM	BMMM
50	12	0.723	6.531	6.296	1.585	6.664	2.360
		(0.041)	(0.194)	(0.147)	(0.079)	(0.246)	(0.131)
100	12	0.717	6.495	6.281	1.590	6.640	2.349
		(0.027)	(0.155)	(0.106)	(0.060)	(0.178)	(0.091)
n	T	$RMS E_H$			RI		
		BCRM	BCM	BMMM	BCRM	BCM	BMMM
50	12	1.631	7.094	6.320	0.931	0.714	0.768
		(0.076)	(0.215)	(0.145)	(0.004)	(0.007)	(0.008)
100	12	1.635	7.068	6.303	0.931	0.715	0.768
		(0.057)	(0.158)	(0.103)	(0.003)	(0.005)	(0.005)

Table 6. Prediction results of each method in C2 configuration.

n	T	$RMS E_L$			$RMS E_U$		
		BCRM	BCM	BMMM	BCRM	BCM	BMMM
50	12	0.738	6.530	6.296	2.913	7.086	3.398
		(0.044)	(0.195)	(0.147)	(0.149)	(0.312)	(0.187)
100	12	0.725	6.494	6.281	2.930	7.072	3.400
		(0.029)	(0.155)	(0.106)	(0.108)	(0.233)	(0.131)
n	T	$RMS E_H$			RI		
		BCRM	BCM	BMMM	BCRM	BCM	BMMM
50	12	2.932	7.749	6.465	0.894	0.708	0.747
		(0.147)	(0.256)	(0.146)	(0.007)	(0.008)	(0.009)
100	12	2.947	7.731	6.450	0.894	0.709	0.747
		(0.108)	(0.190)	(0.102)	(0.005)	(0.006)	(0.006)

Table 7. Prediction results of each method in C3 configuration.

n	T	$RMS E_L$			$RMS E_U$		
		BCRM	BCM	BMMM	BCRM	BCM	BMMM
50	12	1.418	6.645	6.415	1.423	6.625	2.254
		(0.081)	(0.226)	(0.183)	(0.070)	(0.262)	(0.122)
100	12	1.424	6.605	6.394	1.423	6.594	2.240
		(0.055)	(0.176)	(0.129)	(0.056)	(0.181)	(0.084)
n	T	$RMS E_H$			RI		
		BCRM	BCM	BMMM	BCRM	BCM	BMMM
50	12	1.817	7.323	6.440	0.916	0.714	0.770
		(0.073)	(0.231)	(0.181)	(0.005)	(0.007)	(0.008)
100	12	1.821	7.292	6.419	0.916	0.715	0.771
		(0.050)	(0.160)	(0.127)	(0.004)	(0.005)	(0.006)

Table 8. Prediction results of each method in C4 configuration.

n	T	$RMS E_L$			$RMS E_U$		
		BCRM	BCM	BMMM	BCRM	BCM	BMMM
50	12	2.451	6.940	6.723	2.459	6.915	3.015
		(0.140)	(0.267)	(0.231)	(0.120)	(0.305)	(0.152)
100	12	2.464	6.899	6.699	2.462	6.883	3.008
		(0.096)	(0.210)	(0.170)	(0.098)	(0.222)	(0.118)
n	T	$RMS E_H$			RI		
		BCRM	BCM	BMMM	BCRM	BCM	BMMM
50	12	1.817	7.323	6.857	0.860	0.706	0.753
		(0.127)	(0.264)	(0.225)	(0.008)	(0.009)	(0.010)
100	12	1.821	7.292	6.836	0.860	0.707	0.754
		(0.087)	(0.183)	(0.161)	(0.006)	(0.007)	(0.007)

5. Empirical analysis

This section applies the proposed model to the estimation and prediction of the AQI and compares it with other methods. The concentration of all kinds of pollutants changes with space and time. The panel interval-valued data can be used to describe this variation, and this paper aims to construct panel interval-valued data models for the AQI. Based on the AQI-related data in Beijing, Tianjin, Shijiazhuang, and Chongqing in China, this study selected AQI-related data from 4 representative cities for 40 consecutive days (2023.7.20–2023.8.28). Among them, the data of the first 30 days are used to train the models, and the remaining data are used to test the models.

For the panel interval-valued data set $S = (X_{it}, y_{it})$, $i = 1, 2, \dots, N$; $t = 1, 2, \dots, T$, $y_{it} = [y_{it}^L, y_{it}^U]$ is considered to be the observed interval-valued dependent variable, which is the AQI. y_{it}^L represents the minimum value of the AQI in the i th city on date t and y_{it}^U represents the maximum value of the AQI in the i th city on date t . $X_{it} = (X_{it1}, X_{it2}, \dots, X_{it6})^T$ are interval-valued independent vectors, which represent the values of CO , NO_2 , O_3 , PM_{10} , $PM_{2.5}$, and SO_2 , respectively. $X_{itj} = [a_{itj}, b_{itj}]$, $i = 1, 2, 3, 4$, $t = 1, 2, \dots, 40$, $j = 1, 2, \dots, 6$, a_{itj} indicates the minimum value of the j th pollutant in the i th city on date t , b_{itj} indicates the maximum value of the j th pollutant in the i th city on date t .

First of all, the center data and range data of all dependent variables are made into a Q-Q plot and the results are shown in Figure 2. It is considered that the data is approximately normally distributed.

α_i is a random effect and $\beta^s = (\beta_1^s, \beta_2^s, \dots, \beta_6^s)$, where the prior information β_0^s is generated from the linear least squares method. Other prior information is consistent with the numerical simulation. We take the mean value of the last 2000 observations as the estimated value of the parameters based on

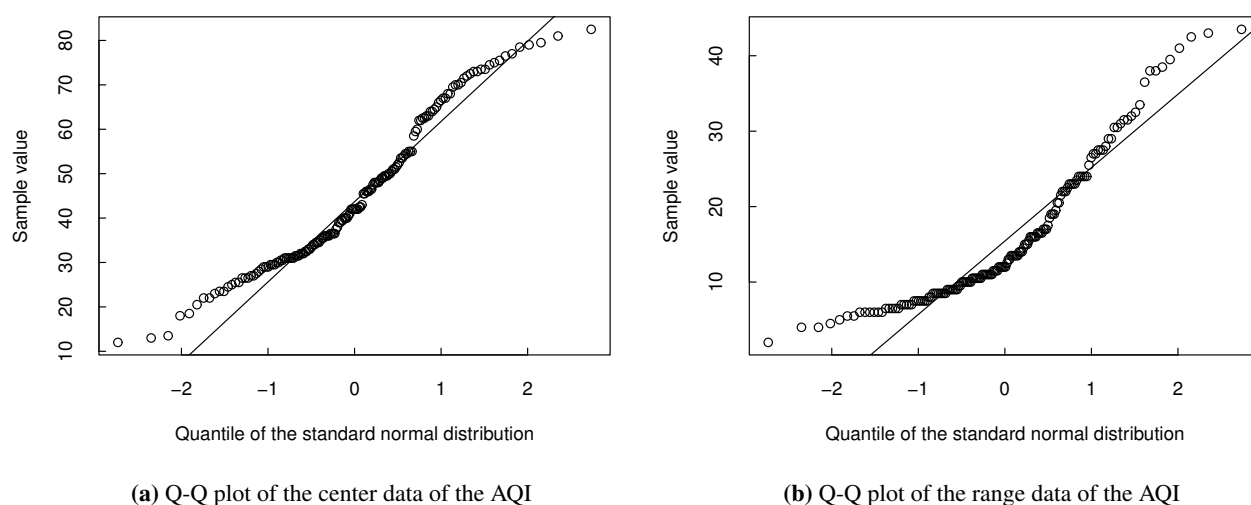


Figure 2. Q-Q plot of the data.

5000 iterations, and the results are shown in Table 9.

Table 9. The results of Bayesian estimation.

Para.	CO^c	NO_2^c	O_3^c	PM_{10}^c	$PM_{2.5}^c$	SO_2^c
Mean	-5.4727	0.1152	0.2858	0.4476	0.2669	-0.5015
SD	2.6396	0.1042	0.0202	0.0820	0.1254	0.2977
Para.	CO^r	NO_2^r	O_3^r	PM_{10}^r	$PM_{2.5}^r$	SO_2^r
Mean	12.5602	-0.5210	0.2283	0.3567	0.2533	-0.1633
SD	2.9272	0.1192	0.0260	0.1236	0.1797	0.4666

It can be seen from the parameters of the central data that NO_2 , O_3 , PM_{10} , and $PM_{2.5}$ have a positive impact on the AQI, that is, the greater the concentration of the pollutants, the greater the AQI and the more severe the air pollution, while CO and SO_2 have a negative impact on the AQI. Similarly, in the range data, CO , O_3 , PM_{10} , and $PM_{2.5}$ have a positive impact on the AQI, while NO_2 and SO_2 have a negative impact on the AQI. The first 75% data (the first 30 days) are selected as the training set, and the last 25% data (the last 10 days) as the test set, and the AQI is predicted by the three methods mentioned above. The prediction results are shown in Table 10.

Table 10. Prediction results of the AQI.

Models	$RMS E_L$	$RMS E_U$	$RMS E_H$	RI
BCRM	8.0494	9.6009	11.5262	0.5820
BCM	12.3272	17.2655	18.7467	0.5109
BMMM	6.8774	9.9299	10.9141	0.6292

From the table, we can see that the BCRM and the BMMM are relatively similar, but consistently better than the BCM. Specifically, the results obtained for the BCRM method are as follows. First, the BCRM method has relatively small values in the indicators of $RMS E_L$, $RMS E_U$, and $RMS E_H$, which proves the accuracy of BCRM method estimation. It reached 0.5820 on RI , which is relatively good. In addition, the BMMM method also predicts well and even better than the BCRM method in some

indicators, which may be due to the fact that the interval-valued data in the empirical study of this paper are presented in the form of minimum and maximum values. Then, comparison charts of different methods based on real data and forecast data are listed in Figure 3, where a solid line represents the predicted data, a dotted line represents the real data, the horizontal coordinate represents the predicted frequency, and the vertical coordinate represents the predicted value of the AQI. As can be seen from Figure 3, the trend of the forecast data based on the BCRM is generally consistent with the real data and the distance between the real data and forecast data is small, which also explains the accuracy of the estimation of the BCRM method to a certain extent. The distance between real data and forecast data is larger for the BCM method, while the distance between real data and forecast data is also smaller for the BMMM method.

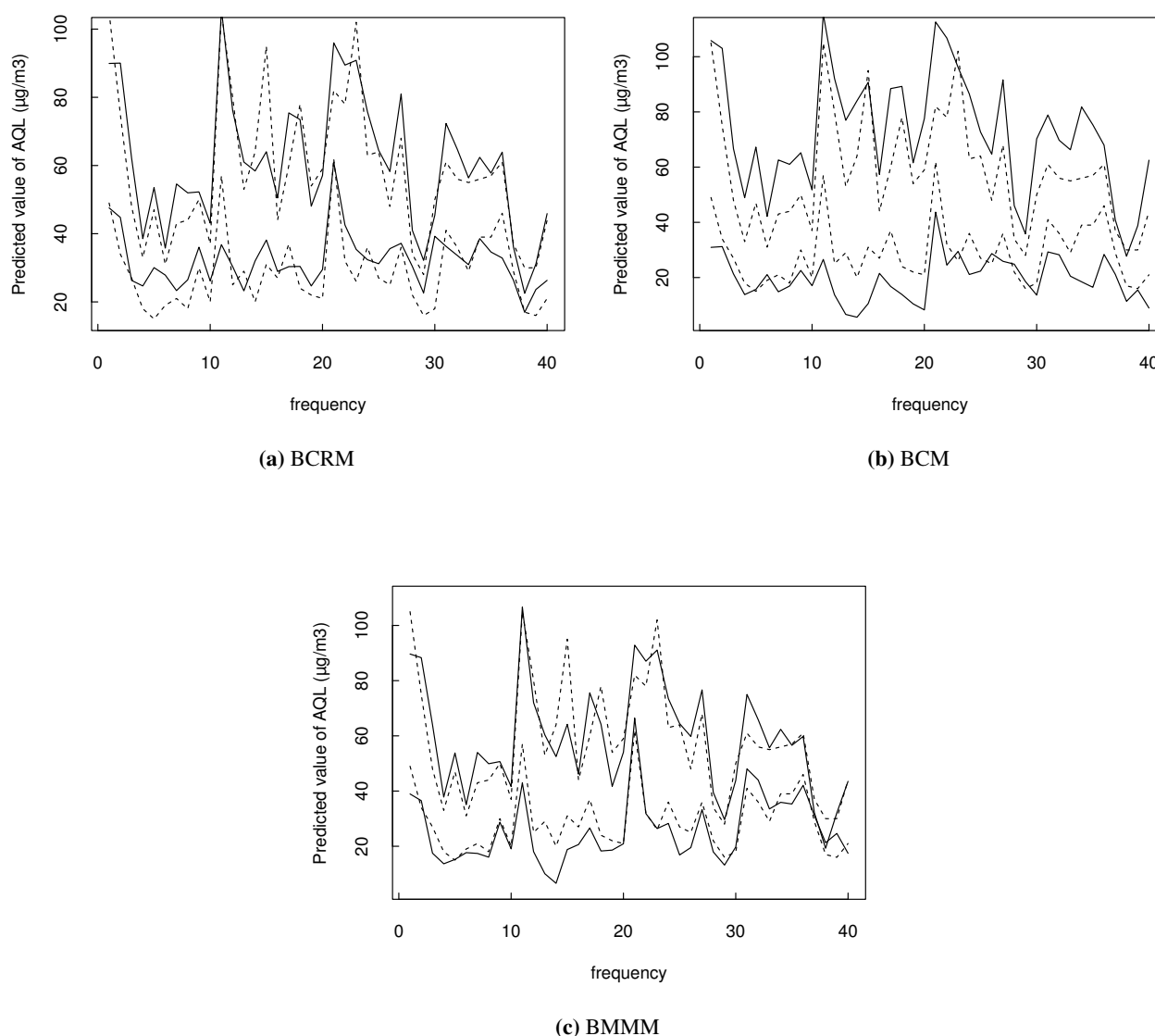


Figure 3. Comparison chart of different methods based on real data and forecast data.

6. Conclusions and discussion

In this paper, a Bayesian model based on the center and range method (BCRM) of random effects panel interval-valued data is proposed and compared with the Bayesian model based on the center method (BCM) of random effects panel interval-valued data and the Bayesian model based on the maximum and minimum method (BMMM) of random effects panel interval-valued data. The results show that: on the one hand, the deviation and standard deviation of parameters estimated by the BCRM model are very small; and on the other hand, in the case of different configurations and different sample sizes, the prediction errors $RMS E_L$, $RMS E_U$, and $RMS E_H$ of the BCRM are basically the smallest, while the interval coverage RI is basically the largest, that is, the BCRM has the best prediction effect. The model is applied to the empirical study of AQI prediction, and the results show the effectiveness of the proposed Bayesian model.

In addition, the model deserves further study. Specific considerations are as follows: Quantile regression has greater flexibility in the distribution of random errors and is therefore robust to non-normal errors and outliers. Therefore, how to combine quantile regression with a Bayesian model of random effect panel interval-valued data is worthy of further study.

Use of AI tools declaration

The authors declare they have not used Artificial Intelligence (AI) tools in the creation of this article.

Acknowledgments

This research was supported by the National Social Science Fund of China (Grant No. 23BTJ069).

Conflict of interest

The authors declare there are no conflicts of interest.

References

1. Y. Y. Sun, X. Y. Zhang, A. T. K. Wan, S. Y. Wang, Model averaging for interval-valued data, *Eur. J. Oper. Res.*, **301** (2022), 772–784. <https://doi.org/10.1016/j.ejor.2021.11.015>
2. L. Billard, E. Diday, *Regression Analysis for Interval-Valued Data, Data Analysis, Classification and Related Methods*, Berlin, Heidelberg: Springer, (2000), 369–374. <https://doi.org/10.1007/978-3-642-59789-3-58>
3. L. Billard, E. Diday, From the statistics of data to the statistics of knowledge: Symbolic data analysis, *J. Am. Stat. Assoc.*, **98** (2003), 470–487. <https://doi.org/10.1198/016214503000242>
4. L. Billard, E. Diday, Descriptive statistics for interval-valued observations in the presence of rules, *Comput. Stat.*, **21** (2006), 187–210. <https://doi.org/10.1007/s00180-006-0259-6>
5. E. D. Lima, F. D. T. de Carvalho, Centre and range method for fitting a linear regression model to symbolic interval data, *Comput. Stat. Data Anal.*, **52** (2008), 1500–1515. <https://doi.org/10.1016/j.csda.2007.04.014>

6. E. D. Lima, F. D. T. de Carvalho, Constrained linear regression models for symbolic interval-valued variables, *Comput. Stat. Data Anal.*, **54** (2010), 333–347. <https://doi.org/10.1016/j.csda.2009.08.010>
7. P. Giordani, Lasso-constrained regression analysis for interval-valued data, *Adv. Data Anal. Classif.*, **9** (2015), 5–19. <https://doi.org/10.1007/s11634-014-0164-8>
8. A. L. S. Maia, F. D. T. de Carvalho, Holt's exponential smoothing and neural network models for forecasting interval-valued time series, *Int. J. Forecasting*, **27** (2011), 740–759. <https://doi.org/10.1016/j.ijforecast.2010.02.012>
9. L. C. Souza, R. M. C. R. Souza, G. J. A. Amaral, T. M. Silva, A parametrized approach for linear-regression of interval data, *Knowl.-Based Syst.*, **131** (2017), 149–159. <https://doi.org/10.1016/j.knosys.2017.06.012>
10. L. T. Kong, X. J. Song, X. M. Wang, Nonparametric regression for interval-valued data based on local linear smoothing approach, *Neurocomputing*, **501** (2022), 834–843. <https://doi.org/10.1016/j.neucom.2022.06.073>
11. L. T. Kong, X. W. Gao, A regularized MM estimate for interval-valued regression, *Expert Syst. Appl.*, **238** (2024). <https://doi.org/10.1016/j.eswa.2023.122044>
12. E. Nuroglu, R. M. Kunst, The effects of exchange rate volatility on international trade flows: evidence from panel data analysis and fuzzy approach, in *Proceedings of Rijeka Faculty of Economics: Journal of Economics and Business*, **30** (2012), 9–31.
13. F. He, D. D. Ma, X. N. Xu, Interval environmental efficiency across provinces in China under the constraint of haze with SBM-undesirable interval model, *J. Arid Land Res. Environ.*, **30** (2016), 28–33. <https://doi.org/10.13581/j.cnki.rdm.20160906.004>
14. S. N. Zhao, R. Q. Liu, Z. F. Shang, Statistical inference on panel data models: A kernel ridge regression method, *J. Bus. Econ. Stat.*, **39** (2019), 325–337. <https://doi.org/10.1080/07350015.2019.1660176>
15. E. Aristodemou, Semiparametric identification in panel data discrete response models, *J. Econom.*, **220** (2021), 253–271. <https://doi.org/10.1016/j.jeconom.2020.04.002>
16. L. R. Liu, H. R. Moon, F. Schorfheide, Forecasting with a panel Tobit model, *Quant. Econom.*, **14** (2023), 117–159. <https://doi.org/10.3982/QE1505>
17. B. H. Beyaztas, S. Bandyopadhyay, Robust estimation for linear panel data models, *Stat. Med.*, **39** (2020), 4421–4438. <https://doi.org/10.1002/sim.8732>
18. H. Liu, Y. Q. Pei, Q. F. Xu, Estimation for varying coefficient panel data model with cross-sectional dependence, *Metrika*, **83** (2020), 377–410. <https://doi.org/10.1007/s00184-019-00739-0>
19. S. Y. Ke, P. C. B. Phillips, L. J. Su, Robust inference of panel data models with interactive fixed effects under long memory: A frequency domain approach, *J. Econom.*, **241** (2024). <https://doi.org/10.1016/j.jeconom.2024.105761>
20. A. B. Ji, J. J. Zhang, X. He, Y. H. Zhang, Fixed effects panel interval-valued data models and applications, *Knowl.-Based Syst.*, **237** (2022). <https://doi.org/10.1016/j.knosys.2021.107798>
21. T. Park, G. Casella. The bayesian lasso, *J. Am. Stat. Assoc.*, **103** (2008), 681–686. <https://doi.org/10.1198/016214508000000337>

22. D. K. Xu, Z. Z. Zhang, A semiparametric Bayesian approach to joint mean and variance model, *Stat. Probab. Lett.*, **83** (2013), 1624–1631. <https://doi.org/10.1016/j.spl.2013.02.023>
23. I. Castillo, J. Schmidt-Hieber, A. V. der Vaart, Bayesian linear regression with sparse priors, *Ann. Stat.*, **43** (2015), 1986–2018. <https://doi.org/10.1214/15-AOS1334>
24. M. Pfarrhofer, P. Piribauer, Flexible shrinkage in high-dimensional Bayesian spatial autoregressive models, *Spatial Stat.*, **29** (2019), 109–128. <https://doi.org/10.1016/j.spasta.2018.10.004>
25. Z. Q. Wang, N. S. Tang, Bayesian quantile regression with mixed discrete and Nonignorable missing covariates, *Bayesian Anal.*, **15** (2020), 579–604. <https://doi.org/10.1214/19-BA1165>
26. D. Zhang, L. C. Wu, K. Y. Ye, M. Wang, Bayesian quantile semiparametric mixed-effects double regression models, *Stat. Theory Relat. Fields*, **5** (2021), 303–315. <https://doi.org/10.1080/24754269.2021.1877961>
27. L. Z. Tang, Y. J. Li, L. J. Zhao, Study on the Bayesian Elastic Net quantile regression for panel data: methods and applications, *Stat. Res.*, **37** (2020), 94–113. <https://doi.org/10.19343/j.cnki.11-1302/c.2020.03.008>
28. C. Q. Tao, Y. T. Xu, Study on Bayesian adaptive lasso quantile regression using asymmetric exponential power distribution for panel data, *Stat. Res.*, **39** (2022), 128–144. <https://doi.org/10.19343/j.cnki.11-1302/c.2022.09.010>
29. J. Zhang, M. Liu, M. Dong, Variational Bayesian inference for interval regression with an asymmetric Laplace distribution, *Neurocomputing*, **323** (2019), 214–230. <https://doi.org/10.1016/j.neucom.2018.09.083>
30. M. Xu, Z. F. Qin, A bivariate Bayesian method for interval-valued regression models, *Knowl.-Based Syst.*, **235** (2022). <https://doi.org/10.1016/j.knosys.2021.107396>
31. J. H. Ding, Z. Q. Zhang, Bayesian Statistical Models with Uncertainty Variables, *J. Intell. Fuzzy Syst.*, **39** (2020), 1109–1117. <https://doi.org/10.3233/JIFS-192014>
32. F. D. T. de Carvalho, R. M. C. R. de Souza, M. Chavent, Y. Lechevallier, Adaptive Hausdorff distances and dynamic clustering of symbolic interval data, *Pattern Recognit. Lett.*, **27** (2006), 167–179. <https://doi.org/10.1016/j.patrec.2005.08.014>
33. C. Y. Hu, L. T. He, An application of interval methods to stock market forecasting, *Reliab. Comput.*, **13** (2007), 423–434. <https://doi.org/10.1007/s11155-007-9039-4>



AIMS Press

©2025 the Author(s), licensee AIMS Press. This is an open access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0>)