

On the Estimation of the Regression Model M for Interval Data

Marta García-Bárzana¹, Ana Colubi¹, and Erricos J. Kontoghiorghes²

Abstract A linear regression model for interval data based on the natural interval-arithmetic has recently been proposed. Interval data can be identified with 2-dimensional points in $\mathbb{R} \times \mathbb{R}^+$, since they can be parametrized by its mid-point and its semi-amplitude or spread, which is non-negative. The model accounts separately for the contribution of the mid-points and the spreads through a single equation. The least squares estimation becomes a quadratic optimization problem subject to linear constraints, which guarantee the existence of the residuals. Several estimators are discussed. Namely, a closed-form estimator, the restricted least-squares estimator, an empirical estimator and an estimator based on separate models for mids and spreads have been investigated. Real-life examples are considered. Simulations are performed in order to assess the consistency and the bias of the estimators. Results indicate that the numerical and the closed-form estimator are appropriate in most of cases, while the empirical estimator and the one based on separate models are not always suitable.

1 Introduction

Often experimental researches involves non-perfect data, as missing data, or censored data. In particular, closed and bounded real-valued sets in \mathbb{R}^p are useful to model information which also representing linguistic descriptions, fluctuations, grouped data images, to name but a few. Interval data are a specific case of this kind of elements. The study of linear regression models working with interval-valued variables has been addressed *mainly* by two ways:

¹ Department of Statistics and Operational Research, University of Oviedo, Spain, garciaarmarta.uo@uniovi.es · colubi@uniovi.es

² Department of Commerce, Finance and Shipping, Cyprus University of Technology, Cyprus, erricos@cut.ac.cy

(a) in terms of the separate models involving some interval components (as the midpoint and the range or the minimum and the maximum) (see Billard and Diday, 2003; Lima Neto *et al.*, 2005 and references therein) which most of the times work with symbolic interval variables; and (b) in terms of arithmetic set-based unified models (as in Diamond 1990, Gil *et al.* 2001, 2002, 2007, González-Rodríguez *et al.* 2007, Blanco-Fernández *et al.* 2011, among others). The main difference between both views is that the first approach usually fits the separate models by numerical or classical tools, but without the usual probabilistic assumptions for the regression model. This provides good fittings but non-obvious easy ways of making inferences. On the other hand, the second approach provides a natural framework to develop inferences, although the least squares approach becomes a minimization problem with strong constraints.

In Blanco-Fernández *et al.* (2011) a flexible simple linear regression model was introduced, the so-called Model M . This model is *flexible* in the sense that it accounts for relationship between mid points and the radius of the involved random intervals. A comparison of several regression estimators of Model M will be addressed.

The rest of the paper is organized as follows: in Section 2 some preliminary about the Model M will be introduced. In Section 3 four estimation approaches of Model M will be described. In Section 4 a real-life example is analyzed to compare the behaviour of the estimators. Finally, Section 5 contains some conclusions.

2 The Model M for Random Intervals

Hereafter, the intervals that will be considered are elements in the space $\mathcal{K}_c(\mathbb{R}) = \{[a_1, a_2] : a_1, a_2 \in \mathbb{R}, a_1 \leq a_2\}$. An interval $A \in \mathcal{K}_c(\mathbb{R})$ can be expressed in terms of its minimum and maximum or in terms of its middle point (mid) and the radius (spr). The second characterization is more usual in regression studies, as it involves non-negativity constraints which are easier to handle than the order constraints involved in the first characterization.

There is another representation for the intervals which will be used, namely, the *canonical decomposition*, defined as $A = \text{mid } A [1 \pm 0] + \text{spr } A [0 \pm 1]$ (see Blanco-Fernández *et al.*, 2011).

The arithmetics which will be used are the *Minkowski addition* $A + B = \{a + b : a \in A, b \in B\}$ and the product by scalars $\lambda A = \{\lambda a : a \in A\}$, with $A, B \in \mathcal{K}_c(\mathbb{R})$ and $\lambda \in \mathbb{R}$. The space $(\mathcal{K}_c(\mathbb{R}), +, \cdot)$ is not linear as the existence of the symmetric element with respect to the addition is not guaranteed in general, in the sense that $A + (-A) \neq \{0\}$ unless A is a singleton. A new concept of difference agreeing with the natural dif-

ference, the so-called *Hukuhara difference*, is introduced. It is defined as $A -_H B = [\inf A - \inf B, \sup A - \sup B]$ if and only if $\text{spr } B \leq \text{spr } A$.

Remark: If $\text{spr } B > \text{spr } A$, then the Hukuhara difference does not exist.

The distance used is the so-called d_τ (see Trutschnig *et al.*, 2009) defined as

$$d_\tau(A, B) = \sqrt{(1 - \tau)(\text{mid } A - \text{mid } B)^2 + \tau(\text{spr } A - \text{spr } B)^2}$$

for all $A, B \in \mathcal{K}_c(\mathbb{R})$.

Random intervals emerged as a generalization of the real-valued random variables. Then, \mathbf{y} is a random interval if it is $\mathcal{B}_{d_\tau} | \mathcal{A}$ measurable, being \mathcal{B}_{d_τ} the Borel σ -algebra and \mathcal{A} the σ -algebra of the probabilistic space (Ω, \mathcal{A}, P) .

Notation: Random intervals will be denoted with boldlowercase letters (\mathbf{x}), vectors with lowercase letters (x) and matrices with uppercase letters (X). The (Aumann) expect value is defined as $E(\mathbf{x}) = [E(\text{mid } \mathbf{x}) \pm E(\text{spr } \mathbf{x})]$, whenever $\text{mid } \mathbf{x}$ and $\text{spr } \mathbf{x} \in L^1(\Omega, \mathcal{A}, P)$. The Aumann expectation fulfils Fréchet principle and the Fréchet variance associated with this expectation is defined as

$$\text{Var}_\tau(\mathbf{x}) = \sigma_{\mathbf{x}, \tau}^2 = E(d_\tau(\mathbf{x}, E(\mathbf{x}))) = (1 - \tau) \sigma_{\text{mid } \mathbf{x}}^2 + \tau \sigma_{\text{spr } \mathbf{x}}^2$$

whenever $\text{mid } \mathbf{x}$ and $\text{spr } \mathbf{x}$ are integrably bounded.

As $(\mathcal{K}_c(\mathbb{R}), +, \cdot)$ is not a linear space, the covariance cannot be defined by mimicking the usual expression involving the arithmetic in $\mathcal{K}_c(\mathbb{R})$. However, it can be defined in \mathbb{R}^2 and we get the following expression

$$\text{Cov}_\tau(\mathbf{x}, \mathbf{y}) = \sigma_{\mathbf{x}, \mathbf{y}} = (1 - \tau) \sigma_{\text{mid } \mathbf{x}, \text{mid } \mathbf{y}} + \tau \sigma_{\text{spr } \mathbf{x}, \text{spr } \mathbf{y}}$$

whenever $\|\text{mid } \mathbf{x}\|_\tau^2, \|\text{mid } \mathbf{y}\|_\tau^2, \|\text{spr } \mathbf{x}\|_\tau^2, \|\text{spr } \mathbf{y}\|_\tau^2 \in L^1(\Omega, \mathcal{A}, P)$.

Model M will relate a response random interval $\mathbf{y} : \Omega \rightarrow \mathcal{K}_c(\mathbb{R})$ with an explanatory random interval $\mathbf{x} : \Omega \rightarrow \mathcal{K}_c(\mathbb{R})$ as follows

$$\mathbf{y} = \mathbf{x}^M \alpha_1 + \mathbf{x}^S \alpha_2 + \varepsilon \tag{1}$$

where $\mathbf{x}^M = \text{mid } \mathbf{x}[1 \pm 0] = [\text{mid } \mathbf{x}, \text{mid } \mathbf{x}]$, $\mathbf{x}^S = \text{spr } \mathbf{x}[0 \pm 1] = [-\text{spr } \mathbf{x}, \text{spr } \mathbf{x}]$, α_1, α_2 and $\varepsilon \in \mathcal{K}_c(\mathbb{R})$ (see Blanco-Fernández *et al.*, 2011).

The Model can be written in the matricial way as

$$\mathbf{y} = \mathbf{x}^{Bl} b_\alpha + \varepsilon \tag{2}$$

with $\mathbf{x}^{Bl} = (\mathbf{x}^M | \mathbf{x}^S) \in \mathcal{K}_c(\mathbb{R})^{1 \times 2}$, $b_\alpha = (\alpha_1 | \alpha_2)^t \in \mathbb{R}^{2 \times 1}$ and $\varepsilon : \Omega \rightarrow \mathcal{K}_c(\mathbb{R})$ being a random interval such that $E(\varepsilon | \mathbf{x}) = \Delta \in \mathcal{K}_c(\mathbb{R})$.

Remark: A property of this model is that it is not identifiable due to the fact that $\mathbf{x}^S = -\mathbf{x}^S$. However, the coefficient α_2 can be considered, without loss of generality, a non-negative vector in \mathbb{R} and the space in which the

solutions to the estimation problem are, can be restricted to \mathbb{R}^+ . In this way, the model is identifiable.

Model M entails the following separate models

$$\begin{aligned}\text{mid } \mathbf{y} &= \alpha_1 (\text{mid } \mathbf{x}) + \text{mid } \varepsilon \\ \text{spr } \mathbf{y} &= |\alpha_2| (\text{spr } \mathbf{x}) + \text{spr } \varepsilon.\end{aligned}\tag{3}$$

Remark: By the assumption that α_2 can be considered non-negative, the second expression can be written as

$$\text{spr } \mathbf{y} = \alpha_2 (\text{spr } \mathbf{x}) + \text{spr } \varepsilon.$$

Thus, it is feasible to consider the estimation of α_1 and α_2 through the estimation of the separate models.

3 Estimation of the Model M

Four estimators of the regression coefficients will be considered. The first one based on the fitting of the separate models introduced in (3). Separate models have already considered to relate interval-valued variables (see Lima Neto & Carvalho 2010 among others). In this case the proposed separate models are:

$$\begin{aligned}\text{mid } \mathbf{y} &= \mathbf{x}^c b^m + \varepsilon^m \\ \text{spr } \mathbf{y} &= \mathbf{x}^s b^s + \varepsilon^s,\end{aligned}\tag{4}$$

where $\mathbf{x}^c = (1, \text{mid } \mathbf{x})$ and $\mathbf{x}^s = (1, \text{spr } \mathbf{x}) \in \mathcal{K}_c(\mathbb{R})^{1 \times 2}$, b^m and $b^s \in \mathbb{R}^{2 \times 1}$, $\mathbf{y} \in \mathcal{K}_c(\mathbb{R})$ and $\varepsilon^m, \varepsilon^s \in \mathbb{R}$. Lima Neto & Carvalho impose the condition that $b^s \geq 0$ to avoid spreads ill-defined. However, b^m has no constraint to be fulfilled.

Then, let $\{\mathbf{y}_i, \mathbf{x}_i\}_{i=1, \dots, n}$ be a random simple sample of intervals, the estimator of b^m will be:

$$\widehat{b}^m = [(x^c)^t (x^c)]^{-1} (x^c)^t \text{mid } \mathbf{y}\tag{6}$$

where $\text{mid } \mathbf{y} \in \mathbb{R}^{n \times 1}$ and

$$x^c = \begin{pmatrix} 1 & 1 & \dots & 1 \\ \text{mid } \mathbf{x}_1 & \text{mid } \mathbf{x}_2 & \dots & \text{mid } \mathbf{x}_n \end{pmatrix} \in \mathbb{R}^{n \times 2}.$$

Parameter b^s is estimated according to Lawson and Hanson algorithm (see Lawson and Hanson, 1974) for constrained LS problems. Then the estimator of both parameters will be denoted by $\widehat{b}_{sep} = (\widehat{b}^m, \widehat{b}^s)$.

Remark: The main drawback of using the separate models to estimate the coefficients is that (5) is not a linear model, due to the non-negativity constraint of the variables. Additionally, the linear independence between

the residuals and the independent variables implies further restrictions on the residuals. Thus, inferences are not straight-forward deduced.

It is possible to obtain another estimator of b_α by using sample moments. Hence, it is introduced the following proposition:

Proposition 1. *Given the random interval \mathbf{y} and the vector of random intervals x^{bl} in the conditions of the Model M , the coefficients' vector b_α can be expressed by:*

$$b_\alpha = Cov_\tau(\mathbf{y}, x^{bl})Cov_\tau(x^{bl}, x^{bl})^{-1}.$$

According to Proposition 1, an empirical estimator could be proposed based on the sample moments, namely:

$$\widehat{b}_{emp} = Cov_\tau(y, X^{bl})Cov_\tau(X^{bl}, X^{bl})^{-1} \quad (7)$$

with $X^{bl} \in \mathcal{K}_c(\mathbb{R})^{2 \times n}$ and $y \in \mathcal{K}_c(\mathbb{R})^{1 \times n}$.

The least squares estimation of b_α and the parameter Δ will be carried out from the information provided by the simple random sample of random intervals $\{\mathbf{y}_i, \mathbf{x}_i\}_{i=1, \dots, n}$ obtained from the model:

$$y = X^{bl} \widehat{b}_\alpha + \widehat{\varepsilon} \quad (8)$$

being

$$X^{bl} = (x^M | x^S) \in \mathcal{K}_c(\mathbb{R})^{n \times 2}$$

and

$$\widehat{b}_\alpha = (\widehat{\alpha}_1 | \widehat{\alpha}_2)^t \in \mathbb{R}^{2 \times 1}.$$

It is necessary to assure the existence of the residuals, or in other words, that the Hukuhara's difference $y -_H (X^{bl} \widehat{b}_\alpha)$ exists. Then the expression of the constraints is:

$$\text{spr}(\widehat{\alpha}_1 x^M + \widehat{\alpha}_2 x^S) \leq \text{spr} y$$

which is equivalent to

$$\text{sign}(\widehat{\alpha}_2) \circ |\alpha_2| \text{spr} x \leq \text{spr} y \equiv \widehat{\alpha}_2 \text{spr} x \leq \text{spr} y.$$

In order to assure the existence of the residuals, the least squares problem will be written as a minimization problem with linear constraints. Specifically, the aim will be to find feasible estimates of b_α and Δ minimizing the not explained variability, that is,

$$\min_{c_2 \in \Gamma} d_\tau^2(y, X^{bl} c + 1\Delta) \quad (9)$$

where $c = (c_1, c_2)^t \in \mathbb{R}^{2 \times 1}$ and $\Gamma = \{c_2 \in [0, \infty) / c_2 \text{spr} x \leq \text{spr} y\}$.

Introducing the following notation, the minimization problem (9) will be transcribed into another one with some useful properties.

$$\begin{aligned} v_m &= \text{mid } y - \overline{\text{mid } \mathbf{y} \mathbf{1}} ; F_m = \text{mid } X^{bl} - \overline{(\text{mid } X^{bl}) \mathbf{1}} \\ v_s &= \text{spr } y - \overline{\text{spr } \mathbf{y} \mathbf{1}} ; F_s = \text{spr } X^{bl} - \overline{(\text{spr } X^{bl}) \mathbf{1}}, \end{aligned} \quad (10)$$

where $v_m, v_s \in \mathbb{R}^{n \times 1}$ and $F_m, F_s \in \mathbb{R}^{n \times 2}$. Then, the minimization problem can be written as:

$$\min_{c_2 \in \Gamma} (1 - \tau)(v_m - F_m c)^t (v_m - F_m c) + \tau (v_s - F_s c)^t (v_s - F_s c). \quad (11)$$

Two possible ways of solving the problem have been proposed. The first one results in a numerical estimator and the second one in an exact expression. Concerning the first approach, as the objective function is a quadratic function and Γ is a set of linear constraints, Karush-Kuhn-Tucker (KKT) Theorem assures the existence of solution and by means of the numerical estimator, which will be denoted in the sequel by \widehat{b}_{kkt} , an estimation of the solution will be obtained.

On the other hand, a closed expression to estimate the regression coefficients has been obtained in Blanco-Fernández *et al.* (2011). It is given in the following proposition and will be the last one to be compared later on.

Proposition 2. *Under the conditions of Model M, the LS regression coefficients estimator is $\widehat{b}_{exact} = (\hat{\alpha}_1, \hat{\alpha}_2)$, where:*

$$\begin{aligned} \hat{\alpha}_1 &= \frac{\text{Cov}(\mathbf{x}^M, \mathbf{y})}{\text{Var}(\mathbf{x}^M)} \\ \hat{\alpha}_2 &= \min \left\{ \hat{a}^{\hat{0}}, \max \left\{ 0, \frac{\text{Cov}(\mathbf{x}^S, \mathbf{y})}{\text{Var}(\mathbf{x}^S)} \right\} \right\} \end{aligned}$$

being $\hat{a}^{\hat{0}} = \min \left\{ \frac{\text{spr } y_i}{\text{spr } x_i} \right\} \forall i \in \{1, \dots, n\}$.

According to Blanco-Fernández *et al.* (2011), given \widehat{b}_α any estimator of b_α , it can be proved that the estimator for the residual, $\widehat{\Delta}$, is:

$$\widehat{\Delta} = \overline{y -_H X^{Bl} \widehat{b}_\alpha},$$

or alternatively as

$$\widehat{\Delta} = \overline{y -_H (x^M \widehat{\alpha}_1 + x^S \widehat{\alpha}_2)}.$$

Indeed, as the existence of Hukuhara's difference $y -_H X^{Bl} \widehat{b}_\alpha$ is guaranteed, $d_\tau^2(y, X^{Bl} \widehat{b}_\alpha + 1\Delta) = d_\tau^2(y -_H X^{Bl} \widehat{b}_\alpha, 1\Delta)$ and applying Fréchet principle, it is obtained

$$\widehat{\Delta} = \overline{y -_H (x^M \widehat{\alpha}_1 + x^S \widehat{\alpha}_2)} = \overline{y -_H (x^M \widehat{\alpha}_1 + x^S \widehat{\alpha}_2)}.$$

4 Applications: A Comparative Study

The first example is concerned with the relationship between the systolic and diastolic pressures in some patients in the hospital Valle del Nalón, in Asturias. The pulse rate as well as both pressure ranges along a day will be modelled by random intervals, where the endpoints of the interval are the minimum and maximum respectively. The mathematical structure will be given by $\Omega = \{3000 \text{ patients of the hospital}\}$, the Borel σ -algebra and a probability P which is uniformly distributed.

Table 1 represents the data of the sample of 56 patients. For this example the constraint $\text{spr } \mathbf{x} b_\alpha \leq \text{spr } \mathbf{y}$ is fulfilled for the 56 patients. Table 2 summarizes the estimates for α_1 and α_2 . For the separate models approach, b_0^m and b_0^s refer to the real-valued intercepts while for the rest of the procedures Δ denotes the interval-valued intercept.

Table 1 \mathbf{y} : diastolic blood pressure (mmHg) and \mathbf{x} : systolic blood pressure (mmHg)

\mathbf{x}	\mathbf{y}	\mathbf{x}	\mathbf{y}	\mathbf{x}	\mathbf{y}
118-173	63-102	119-212	47-93	98-160	47-108
104-161	71-118	122-178	73-105	138-221	70-118
131-186	58-113	127-189	74-125	97-154	60-107
105-157	62-118	113-213	52-112	87-152	50-95
120-179	59-94	141-205	69-133	87-150	47-86
101-194	48-116	99-169	53-109	120-188	53-105
109-174	60-119	126-191	60-98	141-256	77-158
128-210	76-125	99-201	55-121	95-166	54-100
94-145	47-104	88-221	37-94	108-147	62-107
148-201	88-130	94-176	56-121	92-172	45-107
111-192	52-96	102-156	50-94	115-196	65-117
116-201	74-133	103-159	52-95	83-140	45-91
102-167	39-84	102-185	63-118	99-172	42-86
104-161	55-98	111-199	57-113	113-176	57-95
106-167	45-95	130-180	64-121	114-186	46-103
112-162	62-116	103-161	55-97	145-210	100-136
136-201	67-122	125-192	59-101		
90-177	52-104	97-182	54-104		
116-168	58-109	100-161	54-104		
98-157	50-111	159-214	90-127		

All the estimates for α_1 are equal. However, the situation for α_2 is different. Focussing on $\widehat{b_{emp}}$ and $\widehat{b_{exact}}$, it can be seen that they are equal because the sample values fulfil the constraints to assure the existence of the residuals. However, in general, they do not need to be the same value (as shown in next example, due to the fact that $\widehat{b_{exact}}$ was defined to fulfil the constraint, whereas $\widehat{b_{emp}}$ was not). The estimate obtained from the KKT approach is the same as well, but this one was obtained by a numerical approximation. Then we can conclude that the numerical approximation is really close to the exact

Table 2 Estimations of the parameters $\alpha_1, \alpha_2, \Delta$ and b_0^m, b_0^s

Estimator	α_1	α_2	$\Delta/b_0^m - b_0^s$
$\widehat{\mathbf{b}}_{\text{exact}}$	0.4539	0.2570	[1.0164, 32.7000]
$\widehat{\mathbf{b}}_{\text{kkt}}$	0.4539	0.2570	[1.0164, 32.7000]
$\widehat{\mathbf{b}}_{\text{emp}}$	0.4539	0.2570	[1.0164, 32.7000]
$\widehat{\mathbf{b}}_{\text{sep}}$	0.4539	0.6842	16.8582-0.9443

one. \widehat{b}_{sep} reaches a really high value in the estimation of α_2 , which seem to denote that this estimator is not a good one, when it is applied to Model M .

The second example is concerned with the relationship between the familiar average income (\mathbf{y}) and the percentage of people with higher education (\mathbf{x}) in EEUU in 2006 (<http://factfinder.census.gov>). The difference between this example and the previous one is that not all the values of the sample fulfil the constraint $\text{spr } \mathbf{x} b_\alpha \leq \text{spr } \mathbf{y}$. Table 3 displays the data of the sample of 50 people. Then, Table 4 summarizes the values of the different estimates for α_1 and α_2 . Again, estimates of α_1 are equal for all the approaches. However, the estimate of α_2 is different for all the approaches excepting the exact and the KKT-based methods.

5 Conclusions

Some approaches to estimate the regression coefficients have been proposed and the comparison between them have been made by means of some examples. According to the empirical results, estimator \widehat{b}_{sep} does not provide good results, which is natural, as they do not account for the specific features of the unified model that has been considered. Thus, \widehat{b}_{sep} will often divert from the $\widehat{b}_{\text{exact}}$.

The performance of the empirical estimator depends on the data which has been used. If the data satisfies the constraint to assure the existence of the residuals, then the estimator is similar to the exact one. Otherwise, it is an erroneous estimator, as it provides wrong estimates for α_2 , the coefficient accompanying the spreads. In any case, the estimator could be used for large samples, as it approaches to the populational parameter consistently.

Finally, the numerical estimator \widehat{b}_{kkt} is an adequate, as it reaches values which are really close to the exact ones.

Table 3 y : familiar average income, x : percentage of people with higher education

State	y	x	State	y	x
Alabama	48.460-49.954	7.5-7.9	Alaska	67.501-72.243	8.8-10.2
Arizona	55.063-56.355	8.9-9.5	Arkansas	44.28-45.906	5.9-6.5
California	64.150-64.976	10.3-10.5	Colorado	63.639-65.589	12.1-12.7
Connect.	77.203-79.105	14.0-14.8	Delaware	60.406-64.84	9.9-11.1
Columbia	57.076-65.134	24.4-26.4	Florida	54.043-54.847	8.8-9.0
Georgia	55.503-56.721	9.0-9.4	Hawaii	68.823-71.731	9.3-10.3
Idaho	50.612-52.668	6.8-7.4	Illinois	62.592-63.650	10.6-11.0
Indiana	55.322-56.240	7.8-8.2	Iowa	55.158-56.312	7.1-7.7
Kansas	56.159-57.555	9.5-10.1	Kentucky	48.044-49.408	8.0-8.4
Louisiana	47.467-49.055	6.6-7.0	Maine	51.820-53.766	8.5-9.3
Maryland	76.988-78.690	15.4-16.0	Massach.	73.710-75.216	15.4-15.8
Michigan	57.461-58.531	9.0-9.4	Minnesota	66.324-67.294	9.4-9.8
Mississippi	41.797-43.813	5.8-6.4	Missouri	52.465-53.587	8.5-8.9
Montana	50.177-51.835	7.8-9.0	Nebraska	56.291-57.589	8.0-8.8
Nevada	60.629-62.303	6.9-7.5	N.Hampshire	70.065-72.287	10.6-11.8
N.Jersey	77.226-78.524	12.2-12.6	N.Mexico	46.84749.551	10.5-11.3
N.York	61.774-62.502	13.2-13.4	N.Carolina	51.855-52.817	8.1-8.5
N.Dakota	53.918-56.852	5.9-7.1	Ohio	55.760-56.536	8.1-8.5
Oklahoma	47.179-48.731	7.0-7.4	Oregon	55.166-56.680	9.7-10.3
Pennsylv.	57.787-58.509	9.4-9.8	R.Island	62.762-66.704	10.7-11.9
S.Carolina	49.677-50.991	7.7-8.1	S.Dakota	52.870-54.742	6.7-7.7
Tennes.	49.240-50.368	7.3-7.7	Texas	52.080-52.630	7.9-8.1
Utah	57.306-58.976	9.0-9.8	Vermont	56.752-59.574	12.1-13.5
Virginia	66.263-67.509	12.9-13.5	Washington	63.055-64.355	10.5-10.9
W.Virginia	43.189-44.835	6.3-6.9	Wisconsin	60.172-61.096	8.2-8.6
Wyoming	55.797-59.213	6.8-8.0			

Table 4 Estimates of the parameters

Estimator	$\hat{\alpha}_1$	$\hat{\alpha}_2$	$\Delta / \mathbf{b}_0^m - \mathbf{b}_0^s$
$\mathbf{b}_{\text{exact}}$	2.9767	1.3817	[29.7003,30.5269]
\mathbf{b}_{kkt}	2.9767	1.3817	[29.7003,30.5269]
\mathbf{b}_{emp}	2.9767	2.3947	[30.0204,30.2068]
\mathbf{b}_{sep}	2.9767	2.6276	30.1136-0.0196

Acknowledgements This research has been partially supported by the Spanish Ministry of Science and Innovation Grants MTM2009-09440-C02-02 and the Short-Term Scientific Missions associated with the COST Action IC0702 Ref. 010611-008222. Their financial support is gratefully acknowledged. Moreover, part of the work was performed while the first two authors have been visiting the Cyprus University of Technology which partly has supported this research.

References

1. Billard L, Diday E (2003) From the Statistics of data to the Statistics of knowledge: Symbolic Data Analysis. *J Amer Stat Assoc* 98:470–487
2. Blanco-Fernández A, Corral N, González-Rodríguez G (2011) Estimation of a flexible simple linear model for interval data based on set arithmetic. *Comp Stat Data Anal* 55(9):2568–2578
3. Diamond P (1990) Least squares fitting of compact set-valued data. *J Math Anal Appl* 147:531–544
4. Gil MA, Lubiano A, Montenegro M, López-García MT (2002). Least squares fitting of an affine function and strength of association for interval-valued data. *Metrika* 56:97–111
5. González-Rodríguez G, Blanco A, Corral N, Colubi A (2007) Least squares estimation of linear regression models for convex compact random sets. *Adv D Anal Class* 1:67–81
6. Lima Neto EA, DeCarvalho FAT, Freire ES (2005) Applying constrained linear regression models to predict interval-valued data. *LNAI* 3698:92–106
7. Lima Neto EA, DeCarvalho FAT (2010) Constrained linear regression models for symbolic interval-valued variables. *Comp Stat Data Anal* 54:333–347
8. Trutschnig W, González-Rodríguez G, Colubi A, Gil MA (2009) A new family of metrics for compact, convex (fuzzy) sets based on a generalized concept of mid and spread. *Inform Sci* 179(23):3964–3972
9. Lawson CL, Hanson RJ (1974) Solving least squares problems. Prentice-Hall, Englewood Cliffs Reprinted with a detailed “new developments” appendix in 1996 by SIAM Publications, Philadelphia
10. Gil MA, López MT, Lubiano MA, Montenegro M (2001) Regression and correlation analyses of a linear relation between random intervals. *Test* 10(1):183–201 (doi:10.1007/BF02595831)
11. Gil MA, González-Rodríguez G, Colubi A, Montenegro M (2007) Testing linear independence in linear models with interval-valued data. *Computational Statistics and Data Analysis* 51(6):3002–3015. (doi:10.1016/j.csda.2006.01.015)