# On Some Confidence Regions to Estimate a Linear Regression Model for Interval Data

Angela Blanco-Fernández, Norberto Corral,
Gil González-Rodríguez, and Antonio Palacio

**Abstract.** Least-squares estimation of various linear models for interval data has already been considered in the literature. One of these models allows different slopes for mid-points and spreads (or radii) integrated in a unique equation based on interval arithmetic. A preliminary study about the construction of confidence regions for the parameters of that model on the basis of the least-squares estimators is presented. Due to the lack of realistic parametric models for random intervals, bootstrap approaches are proposed. The empirical suitability of the bootstrap confidence sets will be shown by means of some simulation studies.

**Keywords:** Confidence region, Simple linear regression model, Interval random set, Bootstrap approach.

## 1 Introduction

The study of the linear relationship between two random intervals has been addressed in the literature on the basis of several set arithmetic-based regression models (see, for instance, [2, 3, 4, 5, 6, 7, 8]). In order to analyze those models the mid-spread representation of the involved intervals is employed. The utility of this representation is twofold. On one hand, it captures the location and imprecision of the intervals, and on the other hand, it is technically

Angela Blanco-Fernández and Norberto Corral

Statistics and Operational Research Department, Oviedo University,

33007 Oviedo, Spain

e-mail: `blancoangela@uniovi.es,norbert@uniovi.es`

Gil González-Rodríguez and Antonio Palacio

European Centre for Soft Computing, 33600 Mieres, Spain

e-mail: `gil.gonzalez@softcomputing.es,antoniopalacio1982@gmail.com`

easier to handle than the minimum-maximum representation. The linear model presented in [3], denoted by Model M, generalizes those in [4] and [6].

Least squares estimation problems of Model M has been also considered in [3]. On the basis of least-squares estimators different approaches to determine confidence regions can be proposed. Contrary to what happens when the linear regression problem between real random variables is addressed, in the interval scenario no realistic parametric models to describe the distribution of the random sets have been defined up to now. Thus, exact methods are not feasible. Inferential studies about Model M can be developed by means of asymptotic techniques, based on the study of the limit distributions of the regression estimators. To improve the results for finite sample sizes, bootstrap methods are widely considered. In this work several bootstrap approaches are considered in order to build confidence sets for the parameters of the model.

## 2  Preliminaries

Let $(\mathscr{K}_c(\mathbb{R}),+,\cdot)$ be the space of nonempty compact intervals of $\mathbb{R}$ endowed with the semilinear structure induced by the Minkowski addition and the product by a scalar, that is, $A+B = \{a+b \,|\, a \in A, b \in B\}$ and $\lambda A = \{\lambda a \,|\, a \in A\}$ for all $A, B \in \mathscr{K}_c(\mathbb{R})$ and $\lambda \in \mathbb{R}$. Moreover, given $A, B \in \mathscr{K}_c(\mathbb{R})$, if there exists $C \in \mathscr{K}_c(\mathbb{R})$ so that $A = B + C$, then $C$ is defined as the Hukuhara difference between $A$ and $B$, denoted by $A -_H B$. The interval $A$ can be characterized by means of the real vector $(\inf A, \sup A) \in \mathbb{R}^2$ such that $\inf A \leq \sup A$, or equivalently, by its mid-point (or centre) and its spread (or radius), that is, $(\mathrm{mid}A, \mathrm{spr}A)$ with $\mathrm{spr}A \geq 0$, where $\mathrm{mid}A = (\sup A + \inf A)/2$ and $\mathrm{spr}A = (\sup A - \inf A)/2$. The notation $A = [\inf A, \sup A]$ or $A = [\mathrm{mid}A \pm \mathrm{spr}A]$, respectively, will be considered in each case.

Several metrics can be defined on the space $\mathscr{K}_c(\mathbb{R})$. For least squares problems associated with regression studies, an $L_2$-type metric is suitable. Taking inspiration on the family of metrics for compact convex sets introduced in [10], a generalized $L_2$-type distance between two intervals $A$ and $B$ can be defined as

$$d_\theta(A,B) = \sqrt{(\mathrm{mid}A - \mathrm{mid}B)^2 + \theta(\mathrm{spr}A - \mathrm{spr}B)^2} \qquad (1)$$

with $\theta > 0$.

Given a probability space $(\Omega, \mathscr{A}, P)$, a mapping $X : \Omega \to \mathscr{K}_c(\mathbb{R})$ is said to be an *interval-valued random set* (or *random interval*), if it is $\mathscr{A}|\mathscr{B}_{d_\theta}$-measurable, $\mathscr{B}_{d_\theta}$ denoting the $\sigma$-field generated by the topology induced by the metric $d_\theta$ on $\mathscr{K}_c(\mathbb{R})$.

Let $X : \Omega \to \mathscr{K}_c(\mathbb{R})$ be a random interval such that $E(|X|) < \infty$ (with $|X|(\omega) = \sup\{|x| \,|\, x \in X(\omega)\}$ for all $\omega \in \Omega$), then, the *expected value of X in Kudō-Aumann's sense* (see, e.g., [1]) is the interval $E(X) = [E(\inf X), E(\sup X)]$. The variance of $X$ is defined in the classical statistical way, in terms of the $d_\theta$ metric, as $\sigma_X^2 = E\big(d_\theta(X, E(X))^2\big)$, whenever $E(|X|^2) < \infty$. However, it is not possible to define the covariance analogously to the usual concept, due to

the lack of linearity on $\mathscr{K}_c(\mathbb{R})$. Through the *(mid-spr)* parametrization it is possible to define the covariance between $X$ and $Y$ by means of the natural concept of covariance in Hilbert spaces as $\sigma_{X,Y} = E\left(\langle t_X - E_{t_X}, t_Y - E_{t_Y}\rangle_\theta\right)$, where $t_X = (\text{mid}X, \text{spr}X) \in \mathbb{R}^2$ (analogously for $t_Y$), and $\langle\cdot,\cdot\rangle_\theta$ is an inner product on $\mathbb{R}^2$ defined in terms of the constant $\theta > 0$ (see (1)) as $\langle\mathbf{a},\mathbf{b}\rangle_\theta = \mathbf{a}'\begin{pmatrix} 1 & 0 \\ 0 & \theta \end{pmatrix}\mathbf{b}$ for all $\mathbf{a}, \mathbf{b} \in \mathbb{R}^2$. The covariance can be expressed in terms of mids and spreads as $\sigma_{X,Y} = \text{Cov}(\text{mid}X, \text{mid}Y) + \theta\,\text{Cov}(\text{spr}X, \text{spr}Y)$. The variance of the interval $X$ can be also expressed as $\sigma_X^2 = \text{Var}(\text{mid}X) + \theta\,\text{Var}(\text{spr}X)$.

The estimators for the moments of random intervals presented above are the usual ones. Given a simple random sample $\{(X_i, Y_i)\}_{i=1}^n$ from $(X, Y)$, let us define by $\overline{X} = (X_1 + X_2 + \ldots + X_n)/n$, $\widehat{\sigma}_X^2 = \overline{d_\theta(X, \overline{X})^2}$ (analogously $\overline{Y}$ and $\widehat{\sigma}_Y^2$) and $\widehat{\sigma}_{X,Y} = \overline{\langle t_X - \overline{t_X}, t_Y - \overline{t_Y}\rangle_\theta}$ the sample mean, the sample variance and the sample covariance for random intervals, respectively.

## 3   Linear Regression Model M

A natural way to model the relationship between two random intervals has been previously proposed by the expression $Y = \alpha X + \varepsilon$, with $\alpha \in \mathbb{R}$ and $\varepsilon$ a random interval such that $E(\varepsilon|X) = B \in \mathscr{K}_c(\mathbb{R})$ (see [6]). Nevertheless, this model is not flexible enough for many real-life applications. As an example, it can be easily checked that this interval linear model tranfers relationships between the *mid* and *spr* real variables by means of the expressions $\text{mid}Y = \alpha\text{mid}Y + \text{mid}\varepsilon$ and $\text{spr}Y = |\alpha|\text{spr}Y + \text{spr}\varepsilon$. Since both equations involve the same regression coefficient (in absolute value), the model is somehow restrictive.

With the aim of considering the *mid* and *spr* components of the intervals separately, but keeping the good properties of the interval arithmetic, a new representation has been introduced in [3]. Each interval $A \in \mathscr{K}_c(\mathbb{R})$ can be expressed as $A = \text{mid}A[1 \pm 0] + \text{spr}A[0 \pm 1]$. This notation gives the inspiration to formalize the called **Model M** between $X$ and $Y$ in [3] as

$$Y = \alpha\,\text{mid}X[1 \pm 0] + \beta\,\text{spr}X[0 \pm 1] + \varepsilon \tag{2}$$

with $\alpha, \beta \in \mathbb{R}$ and $E(\varepsilon|X) = B \in \mathscr{K}_c(\mathbb{R})$. For simpler notation, the linear model (2) will be denoted by $Y = \alpha X^M + \beta X^S + \varepsilon$. Moreover, it is easy to check that $X^S = -X^S$ (since $X^S(\omega) = [-\text{spr}X(\omega), \text{spr}X(\omega)]$, for all $\omega \in \Omega$), so it is possible to consider $\beta \geq 0$ without loss of generality.

From (2) the linear relationships for *mid* and *spr* variables of $X$ and $Y$ are $\text{mid}Y = \alpha\text{mid}Y + \text{mid}\varepsilon$ and $\text{spr}Y = |\beta|\text{spr}Y + \text{spr}\varepsilon$, which clearly entails more flexibility. The flexibility is associated with the extra parameter of Model M, which depends on two scalars and one interval value.

### 3.1   Least Squares Estimation of Model M

The least-squares (LS) estimation of the regression parameters of the model
(2) has been developed in [3]. The LS approach leads to a contrained mini-
minization problem, namely,

$$\left.(\widehat{\alpha},\widehat{\beta},\widehat{B}) = \operatorname*{argmin}_{\left\{a\in\mathbb{R},b\geq 0,C\in\mathscr{K}_c(\mathbb{R})\right\}} \frac{1}{n}\sum_{i=1}^{n} d_\theta^2(Y_i, aX_i^M + bX_i^S + C)\atop \text{subject to}\ \ b\in S\right\} \quad (3)$$

where $S = \{b\in[0,\infty) : Y_i -_H bX_i^S \text{exists, for all } i = 1,\ldots,n\}$. It is easy to check
that $b\in S$ implies that $Y_i -_H (aX_i^M + bX_i^S)$ exists for all $i = 1,\ldots,n$ and for all
$a\in\mathbb{R}$. The existence of these Hukuhara differences assures the existence of
the residuals of the sample model, and thus, the coherence of the solutions
as suitable estimators of the regression parameters.

It should be underlined that, as it was shown in [6] for the simpler model,
if the restriction is overseen, the obtained estimates of the parameters could
not work as estimates for the model (because the residuals could not exist).

The resolution of problem (3) provides the following expressions:

$$\widehat{\alpha} = \frac{\widehat{\sigma}_{X^M,Y}}{\widehat{\sigma}_{X^M}^2}\ ,\quad \widehat{\beta} = \min\left\{\widehat{s}_0, \max\left\{0, \frac{\widehat{\sigma}_{X^S,Y}}{\widehat{\sigma}_{X^S}^2}\right\}\right\}\ \text{and} \quad (4)$$

$$\widehat{B} = \overline{Y} -_H \left(\widehat{\alpha}\overline{X^M} + \widehat{\beta}\overline{X^S}\right)\ ,$$

where $\widehat{s}_0 = \min\{\operatorname{spr}Y_i/\operatorname{spr}X_i : \operatorname{spr}X_i \neq 0\}$ ($\widehat{s}_0 = \infty$ if $\operatorname{spr}X_i = 0$ for all $i = 1,\ldots,n$).

## 4   Bootstrap Confidence Regions for the Regression Parameters

Since it is not feasible to look for the exact distribution of the LS estimators
and since the asymptotic results usually provide good results only for very
large sample sizes, in this section some alternatives based on bootstrapping
are explored.

Different schemes to generate bootstrap samples from Model (2) can be
followed. When a fixed design is considered (that is, the independent variable
is not random but deterministic), the most usual procedure is the *residual
bootstrap*. On other hand, when both variables in the linear model are con-
sidered as random elements, the natural resampling is made from a simple
random sample of the pair of variables by means of the *paired bootstrap* (see
[9] for a complete description of both procedures). The linear model (2) is
formalized for two random intervals, so the paired bootstrap approach will
be used for the development of inferential studies about Model M.

Several bootstrap confidence sets can be constructed for the regression
parameters of linear models involving real-valued random variables (see [9]).

The best known ones are the *percentile*, *hybrid* and *t-* bootstrap confidence set. Each of them is based on the sample distribution of a different bootstrap expression obtained from the bootstrap estimator of the parameter.

Let $X$ and $Y$ be random intervals verifying Model (2). The separate expressions for the least-squares estimators of the parameters $\alpha$ and $\beta$ presented in (4) allow us to build confidence sets for each parameter separately. Let $\widehat{\alpha}$ be the least-squares estimator of $\alpha$ obtained from a simple random sample $\{X_i, Y_i\}_{i=1}^{n}$ from $(X, Y)$. We denote by $\{X_i^*, Y_i^*\}_{i=1}^{n}$ a bootstrap sample, generated by means of the election of $n$ elements uniformly and with replacement from $\{X_i, Y_i\}_{i=1}^{n}$. Let $\widehat{\alpha}^*$ be the least-squares estimator of $\alpha$ with respect to the bootstrap sample. From the bootstrap estimator $\widehat{\alpha}^*$ the procedure to build the three confidence intervals (CI) for parameter $\alpha$ follows.

- *Bootstrap percentile CI*: If we denote by $K_{BOOT}$ the distribution function of the bootstrap estimator $\widehat{\alpha}^*$, the bootstrap percentile CI for $\alpha$ at a confidence level $1 - \rho$ is defined by means of the corresponding percentiles of $K_{BOOT}$, that is,

$$IC_P(\alpha)_{1-\rho} = \left[ \; K_{BOOT}^{-1}(\rho/2) \; , \; K_{BOOT}^{-1}(1 - \rho/2) \; \right], \qquad (5)$$

  where $K_{BOOT}^{-1}$ denotes the pseudoinverse of $K_{BOOT}$.

- *Bootstrap hybrid CI*: Let $H_{BOOT}$ be the distribution function of the term $n^l(\widehat{\alpha}^* - \widehat{\alpha})$, where $l$ is an arbitrary constant. $H_{BOOT}(x) = P[n^l(\widehat{\alpha}^* - \widehat{\alpha}) \leq x]$, for $x \in \mathbb{R}$. The most usual election for $l$ is $1/2$. Thus, the bootstrap hybrid CI for $\alpha$ at significance level $\rho$ has got the expression

$$IC_H(\alpha)_{1-\rho} = \left[ \; \widehat{\alpha} - \frac{1}{\sqrt{n}} H_{BOOT}^{-1}(1 - \rho/2) \; , \; \widehat{\alpha} - \frac{1}{\sqrt{n}} H_{BOOT}^{-1}(\rho/2) \; \right] \qquad (6)$$

- *t-bootstrap CI*: We consider the standarized pivot $R = \dfrac{\widehat{\alpha} - \alpha}{\widehat{\sigma}_{\widehat{\alpha}}}$ , where $\widehat{\sigma}_{\widehat{\alpha}}^2$ is an estimator of the variance of $\widehat{\alpha}$, and the bootstrap replica of $R$, $R^* = \dfrac{\widehat{\alpha}^* - \widehat{\alpha}}{\widehat{\sigma}_{\widehat{\alpha}^*}^*}$ , with $\widehat{\sigma}_{\widehat{\alpha}^*}^*$ the analogous estimator for the variance of $\widehat{\alpha}^*$. If we denote by $G_{BOOT}$ the distribution function of $R^*$, the t-bootstrap CI for $\alpha$ at confidence level $1 - \rho$ is given by

$$IC_T(\alpha)_{1-\rho} = \left[ \; \widehat{\alpha} - \widehat{\sigma}_{\widehat{\alpha}} G_{BOOT}^{-1}(1 - \rho/2) \; , \; \widehat{\alpha} - \widehat{\sigma}_{\widehat{\alpha}} G_{BOOT}^{-1}(\rho/2) \; \right] \qquad (7)$$

*Remark 1.* The percentiles of the functions $K_{BOOT}$, $H_{BOOT}$ and $G_{BOOT}$ (in each case) can be approximated from the empirical distribution of $\widehat{\alpha}^*$ by means of MonteCarlo Method.

*Remark 2.* It can be shown that the estimator of the variance of $\widehat{\alpha}$ can be expressed as $\widehat{\sigma}_{\widehat{\alpha}}^2 = \dfrac{\widehat{\sigma}_{\text{mid}\widehat{\varepsilon}}^2}{n\widehat{\sigma}_{\text{mid}X}^2}$. However, it is difficult to obtain an analytic expression for $\widehat{\sigma}_{\widehat{\beta}}$. In this case, a bootstrap estimator of the variance of $\widehat{\beta}$ can be approximated by means of MonteCarlo Method based on $B_2$ bootstrap replications (see [9]).

Taking into account the definitions and remarks presented above, an algorithm for the construction of the *percentile*, *hybrid* and *t*- bootstrap confidence set for parameter $\alpha$ of Model M has the following form.

### Algorithm: bootstrap confidence sets for $\alpha$

Let $\{X_i, Y_i\}_{i=1}^n$ be a random sample obtained from $(X, Y)$. Let $\rho$ be a fixed significance level and $B \in \mathbb{N}$ large enough.

P1. Compute the estimates $\widehat{\alpha}$ and $\widehat{\sigma}_{\widehat{\alpha}}^2$.

P2. Generate $B$ bootstrap samples $\{X_i^*, Y_i^*\}_{i=1}^n$ of size $n$, resampling with replacement from the original sample $\{X_i, Y_i\}_{i=1}^n$.

P3. For each iteration $b = 1, \ldots, B$, compute the estimate for $\alpha$ from the corresponding bootstrap sample, $\widehat{\alpha}^*(b) = \dfrac{\widehat{\sigma}_{X^{M^*}, Y^*}}{\widehat{\sigma}_{X^{M^*}}^2}$, and the bootstrap estimator of its variance, $\widehat{\sigma}_{\widehat{\alpha}^*}^{*2} = \dfrac{\widehat{\sigma}_{\text{mid}\varepsilon^*}^2}{n\widehat{\sigma}_{\text{mid}X^*}^2}$.

P4. Aproximate the lower and upper limits of the intervals (5), (6) and (7) substituting the quantiles of the distributions with the corresponding quantiles from the empirical distribution of $\widehat{\alpha}^*$. That is, the values $\{\widehat{\alpha}^*(b)\}_{b=1}^B$ are increasing ordered, and the ones in position $[(\rho/2)B]+1$ and $[(1-\rho/2)B]$ are selected (where $[\cdot]$ denotes the integer function). Let $\widehat{\alpha}_{C1}^*$ and $\widehat{\alpha}_{C2}^*$ be that values. Thus, the *percentile*, *hybrid* and *t*- confidence sets for $\alpha$ at a confidence level $1-\rho$ are given by

$$IC_P(\alpha)_{1-\rho} = \left[\ \widehat{\alpha}_{C1}^*\ ,\ \widehat{\alpha}_{C2}^*\ \right],$$

$$IC_H(\alpha)_{1-\rho} = \left[\ 2\widehat{\alpha} - \widehat{\alpha}_{C2}^*\ ,\ 2\widehat{\alpha} - \widehat{\alpha}_{C1}^*\ \right], \text{ and}$$

$$IC_T(\alpha)_{1-\rho} = \left[\ \widehat{\alpha} - \widehat{\sigma}_{\widehat{\alpha}}\frac{\widehat{\alpha}_{C2}^* - \widehat{\alpha}}{\widehat{\sigma}_{\widehat{\alpha}_{C2}^*}^*}\ ,\ \widehat{\alpha} - \widehat{\sigma}_{\widehat{\alpha}}\frac{\widehat{\alpha}_{C1}^* - \widehat{\alpha}}{\widehat{\sigma}_{\widehat{\alpha}_{C1}^*}^*}\ \right]$$

respectively.

An analogous algorithm can be developed for the construction of the bootstrap confidence sets for the regression parameter $\beta$ in Model M, taking into account the details explained in Remark 2.

## 4.1 Simulation Studies

The empirical behaviour of the bootstrap procedure can be shown by means of some simulation studies. Let us define a theoretical situation for two random intervals $X$ and $Y$ associated by means of the expression

$$Y = X^M + X^S + \varepsilon \tag{8}$$

where the independent interval $X$ is characterized through the real random vector $(\mathrm{mid}X, \mathrm{spr}X)$ such that $\mathrm{mid}X \sim N(0,1)$ and $\mathrm{spr}X \sim \chi_1^2$, and the error interval term is also defined by $\mathrm{mid}\varepsilon \sim N(0,1)$ and $\mathrm{spr}\varepsilon \sim \chi_1^2 + 1$ independent from $X$.

For different samples sizes $n$, a random sample from $(X,Y)$ is simulated. Let $\{X_i, Y_i\}_{i=1}^n$ be one of them. For $k = 10000$ iterations of the suggested bootstrap algorithms, the $0.95 -$ bootstrap confidence sets for $\alpha$ (and analogously for $\beta$) based on $B = 1000$ bootstrap replications are computed, checking for each of them if the theoretical parameter $\alpha = 1$ (and $\beta = 1$) belongs to the corresponding confidence interval. Finally, the coverage rates are gathered in Table 1.

**Table 1** Empirical confidence level of the bootstrap CIs for $\alpha$ and $\beta$

|  | $\alpha$ | | | $\beta$ | | |
| n | $IC_P(\alpha)$ | $IC_H(\alpha)$ | $IC_t(\alpha)$ | $IC_P(\beta)$ | $IC_H(\beta)$ | $IC_t(\beta)$ |
| 30 | 0.9301 | 0.9318 | 0.9374 | 0.8852 | 0.8911 | 0.8969 |
| 50 | 0.9360 | 0.9458 | 0.9466 | 0.8985 | 0.9061 | 0.9067 |
| 100 | 0.9460 | 0.9465 | 0.9476 | 0.9012 | 0.9082 | 0.9124 |
| 200 | 0.9475 | 0.9487 | 0.9494 | 0.9111 | 0.9123 | 0.9152 |

Since the success rates are close to the nominal confidence level $0.95$ (the larger sample size, the closer they are), the empirical correctness of the bootstrap procedure is justified. Indeed, for parameter $\alpha$, the rate of convergence of the empirical significance level can be found in [9]. $IC_t(\alpha)$ is the most accurate, $IC_H(\alpha)$ the second one, and $IC_P(\alpha)$ is the less accurate of the three approaches. In the case of parameter $\beta$, the approximation to the nominal level is slower. A preliminary analysis of this result has shown that the expression of the estimator $\widehat{\beta}$ depending on the sample term $\widehat{s}_0$ entails that the bootstrap estimator $\widehat{\beta}^*$ does not always perform well. Let us recall that $\widehat{s}_0$ is an order statistic (it is defined as the minimum of several real random variables), for which classic bootstrap methods are inconsistent in some situations (see [9]).

## 5   Concluding Remarks

Different procedures to construct bootstrap confidence sets for the parameters $\alpha$ and $\beta$ of Model M have been proposed. Their empirical correctness has been shown by means of some simulation studies. With respect to the parameter $\beta$, a wider study and a possible improvement of the bootstrap procedure for the construction of confidence sets will be addressed in future research. The statistical study of Model M will be extended by means of the development of other inferential studies, like hypothesis testing, the study of linear independence, among others.

## References

1. Aumann, R.J.: Integrals of set-valued functions. J. Math. Anal. Appl. 12, 1–12 (1965)
2. Blanco-Fernández, A., Colubi, A., Corral, N., González-Rodríguez, G.: On a linear independence test for interval-valued random sets. In: Dubois, D., Lubiano, M.A., Prade, H., Gil, M.A., Grzegorzewski, P., Hryniewicz, O. (eds.) Soft Methods for Handling Variability and Imprecision. Advances in Soft Computing, vol. 48, pp. 331–337. Springer, Heidelberg (2008)
3. Blanco-Fernández, A., Corral, N., González-Rodríguez, G.: Estimation of a flexible simple linear model for interval data based on the set arithmetic (submitted for publication, 2010)
4. Gil, M.A., Lubiano, M.A., Montenegro, M., López-García, M.T.: Least squares fitting of an affine function and strength of association for interval-valued data. Metrika 56, 97–111 (2002)
5. Gil, M.A., González-Rodríguez, G., Colubi, A., Montenegro, M.: Testing linear independence in linear models with interval-valued data. Comput. Statist. Data Anal. 51(6), 3002–3015 (2007)
6. González-Rodríguez, G., Blanco-Fernández, A., Corral, N., Colubi, A.: Least squares estimation of linear regression models for convex compact random sets. Adv. Data Anal. Class. 1, 67–81 (2007)
7. Montenegro, M., Casals, M.R., Lubiano, M.A., Gil, M.A.: Two-sample hypothesis tests of means of a fuzzy random variable. Inf. Sci. 133, 89–100 (2001)
8. Montenegro, M., Colubi, A., Casals, M.R., Gil, M.A.: Asymptotic and Bootstrap techniques for testing the expected value of a fuzzy random variable. Metrika 59(1), 31–49 (2004)
9. Shao, J., Tu, D.: The Jackknife and Bootstrap. Springer, New York (1995)
10. Trutschnig, W., González-Rodríguez, G., Colubi, A., Gil, M.A.: A new family of metrics for compact convex (fuzzy) sets based on a generalized concept of mid and spread. Inf. Sci. 179(23), 3964–3972 (2009)