

## GENERALIZATION OF JACCARD INDEX FOR INTERVAL DATA ANALYSIS

A. N. Bazhenov<sup>1,2</sup> and A. Yu. Telnova<sup>1</sup>

UDC 519.612

*Data samples with interval uncertainty are analyzed. It is proposed to use the Jaccard measure (index), which is widely used when comparing sets in various problem areas, as a measure (functional) of the consistency of interval values and their samples. Information about interval analysis, classical and complete (Kaucher) interval arithmetic is presented. For interval quantities, the necessary concepts and definitions of operations are introduced, in particular, generalizations of the concepts of intersection and union of sets. The Jaccard measure is generalized to the case of data with interval uncertainty and samples of interval data. Various variants of interval relations are described in detail — from their coincidence to incompatible cases. Various definitions of the Jaccard measure are given, both symmetric and nonsymmetric with respect to the operands. The connections of the proposed measure with the interval mode and the results of calculations with tweens are considered. A practical example of finding the information set of an interval problem using a new measure is given. Two areas of application of both symmetric and asymmetric measures are presented — computational processes (for characterizing iterative computational processes) and data analysis (for characterizing measurement workspaces and classifying data by a set of features).*

**Keywords:** Jaccard coefficient,  $t$ -norm, consistency measure, interval analysis, covering and noncovering measurements and samples, information set, interval mode, interval tweens.

**Introduction.** One of the tasks of data analysis is to find out the degree of similarity of the sample elements among themselves or the degree of similarity of different samples. To characterize the degree of similarity, it is necessary to introduce a quantitative measure (similarity coefficient). Let us define the measure of similarity formally. The measure of similarity is a two-place function  $S(A, B) \rightarrow [0, 1]$  — a real function on the set of pairs of objects  $(A, B)$ , where  $A, B$  are arbitrary finite sets. Formally, membership of an operation to similarity measures is determined by a system with the following properties:

- boundedness  $0 \leq S(A, B) \leq 1$ ;
- symmetry  $S(A, B) = S(B, A) \geq 0$ ;
- reflexivity  $S(A, A) = 1 \Leftrightarrow A = A$ ;
- transitivity  $A \subseteq B \subseteq C \Rightarrow S(A, B) \geq S(A, C)$ .

There are other systems of similarity axioms. The listed properties, also called the  $t$ -norm, can be implemented in a variety of ways. Measures of inclusion and similarity have been widely used in practice since the second half of the last century. In the publication [1], various similarity measures and their mutual relationships are considered in detail. The first such measure was proposed by Jaccard in 1901 in relation to biology. The Jaccard coefficient  $J$  (index) for sets  $A, B$  has the form

$$J(A, B) = \frac{n(A \cap B)}{n(A \cup B)}, \quad (1)$$

where  $n$  is the cardinality (number of representatives) of the sets.

<sup>1</sup>Loffe Institute, St. Petersburg, Russia; e-mail: bazhenov\_an@spbstu.ru; anna.telnova@mail.ioffe.ru.

<sup>2</sup>Peter the Great St. Petersburg Polytechnic University, Institute of Applied Mathematics and Mechanics, St. Petersburg, Russia.

In computer applications (image processing, machine learning), the measure of similarity of sets is often referred to as IoU (intersection over union).

The purpose of this work is to give a mathematical definition of the proposed option for calculating the binary measure of similarity of sets, to consider various options for its calculation and to study its properties. The measure is generalized to work with data samples with interval uncertainty. The article gives a number of examples and presents various directions for further research.

**Interval arithmetic.** In this publication, similarity measures are applicable to mathematical objects with interval uncertainty. The modern notation system for interval objects used in this article is given in [2]. Here and throughout the article, interval objects are indicated in bold italic type, as is customary in interval arithmetic.

Classical interval arithmetic  $\mathbb{IR}$  is an algebraic system formed by intervals  $x = [\underline{x}, \bar{x}] \subset \mathbb{R}$  so that for any arithmetic operation "\*" from the set  $\{+, -, \cdot, /\}$  the result of the operation between intervals is defined as  $\mathbf{x}*\mathbf{y} = \{x*y|x \in \mathbf{x}, y \in \mathbf{y}\}$ , where  $x, y$  are real numbers;  $\mathbf{x}, \mathbf{y}$  are intervals. The symbol  $\mathbb{I}$  in the notation  $\mathbb{IR}$  is used to distinguish interval and real arithmetic, with  $\mathbb{IR} \supseteq \mathbb{R}$ .

We introduce the characteristics of the interval  $\mathbf{a}$ :

$$\begin{aligned} \text{middle mid } \mathbf{a} &= (\bar{\mathbf{a}} + \underline{\mathbf{a}})/2; \\ \text{radius rad } \mathbf{a} &= (\bar{\mathbf{a}} - \underline{\mathbf{a}})/2; \\ \text{width wid } \mathbf{a} &= \bar{\mathbf{a}} - \underline{\mathbf{a}}, \end{aligned} \tag{2}$$

where  $\bar{\mathbf{a}}, \underline{\mathbf{a}}$  are the right and left (or upper and lower) boundaries of the interval  $\mathbf{a}$ .

In a more general setting, the problem of interval analysis can be solved in the so-called complete interval arithmetic or Kaucher arithmetic  $\mathbb{KR}$  [3]. This arithmetic is obtained by adding improper intervals  $[\underline{x}, \bar{x}]$ ,  $\underline{x} > \bar{x}$  to the set  $\mathbb{IR}$ . The introduction of Kaucher's interval arithmetic is motivated by a number of limitations of classical interval arithmetic: the absence of inverse elements relative to addition and multiplication operations; the impossibility in the general case of taking the minimum by inclusion; the complexity with the formulation of minimax problems. The use of improper intervals allows you to build original methods for solving problems with interval objects.

**Generalization of the Jaccard coefficient.** To analyze data, it is necessary to compare interval objects in a universal way, regardless of the property of the covering. Let us introduce a numerical characteristic of the degree of consistency of two intervals  $\mathbf{x}, \mathbf{y} \subset \mathbb{IR}$  in the form  $JK(\mathbf{x}, \mathbf{y})$  ( $K$  indicates the relation to Kaucher's interval arithmetic):

$$JK(\mathbf{x}, \mathbf{y}) = \frac{\text{wid}(\mathbf{x} \wedge \mathbf{y})}{\text{wid}(\mathbf{x} \vee \mathbf{y})}, \tag{3}$$

where the measure for the Jaccard coefficient, in contrast to (1), is the width of the interval (2); and instead of the operations of intersection and union of sets, we have the operations of taking the minimum  $\wedge$  and maximum  $\vee$ , respectively, by including two quantities in the complete interval arithmetic (Kaucher), described by

$$\begin{aligned} \mathbf{a} \wedge \mathbf{b} &= [\max\{\underline{\mathbf{a}}, \underline{\mathbf{b}}\}, \min\{\bar{\mathbf{a}}, \bar{\mathbf{b}}\}]; \\ \mathbf{a} \vee \mathbf{b} &= [\min\{\underline{\mathbf{a}}, \underline{\mathbf{b}}\}, \max\{\bar{\mathbf{a}}, \bar{\mathbf{b}}\}]. \end{aligned}$$

In the general case, the inclusion minimum in expression (3) may be an improper interval of negative width. To further generalize measure (3) to sets of intervals, we write its numerical expression

$$JK(\mathbf{x}, \mathbf{y}) = \frac{\min\{\bar{\mathbf{x}}, \bar{\mathbf{y}}\} - \max\{\underline{\mathbf{x}}, \underline{\mathbf{y}}\}}{\max\{\bar{\mathbf{x}}, \bar{\mathbf{y}}\} - \min\{\underline{\mathbf{x}}, \underline{\mathbf{y}}\}}. \tag{4}$$

In writing formula (4), instead of interval widths, we use explicit expressions for taking the minimum and maximum by inclusion, which provide the universal character of (4), regardless of whether the result of the operation of taking the minimum by inclusion  $\wedge$  is a proper or improper interval.

Let us illustrate various measures of binary consistency. Consider different cases of the relative position of intervals  $x, y$  and various measures of consistency. Suppose

$$\min \{x, y\} = 1; \quad \max \{\bar{x}, \bar{y}\} = 6 .$$

In Fig. 1 we show the possible options for the relative position of the intervals  $x, y$  and the corresponding values of the measure  $JK$  (cf. formula (4)). Summing up the indicated options, we obtain the inequality

$$-1 \leq JK(x, y) \leq 1 .$$

For non-empty intersections of intervals, the classical Jaccard measure (1) (not shown in Fig. 1) and the new measure  $JK(x, y)$  (3), (4) give the same results. For empty intersections of intervals, the classical Jaccard measure is equal to zero, and  $JK(x, y)$  numerically describes the measure of inconsistency from 0 to  $-1$ .

**Asymmetric variants of the measure  $JK(x, y)$ .** Measure (3) is symmetric with respect to its arguments, which may be inconvenient in case of a large difference in the widths of the arguments. It is desirable that, in this case, the sensitivity of the measure also differ significantly when it is taken according to different arguments. In this case, asymmetric options are useful:

$$JK_x(x, y) = \frac{\text{wid}(x \wedge y)}{\text{wid } x}, \tag{5}$$

$$JK_y(x, y) = \frac{\text{wid}(x \wedge y)}{\text{wid } y}. \tag{6}$$

Similar expressions in the literature are called similarity measures (similarity, overlapping ratio [4, 5]). Expressions (5), (6) differ from similar ones given in [4, 5] in that they allow the interval to be improper when taking the minimum by inclusion.

*Example.* Let  $x = [1; 2], y = [3; 7]$ . Calculate the measures (5), (6):

$$JK_x(x, y) = \frac{2 - 3}{1} = -1, \tag{7}$$

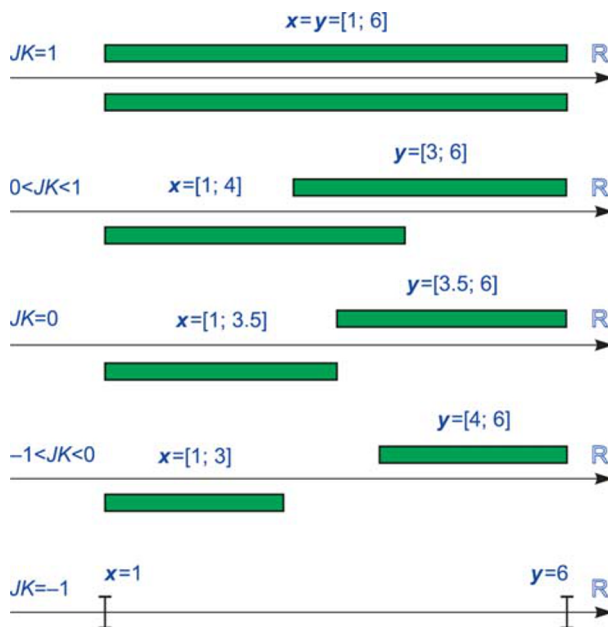


Fig. 1. Variants of relations between the positions of two intervals and the value of the Jaccard measure  $JK$ .

$$JK_y(\mathbf{x}, \mathbf{y}) = \frac{2 - 3}{4} = \frac{-1}{4}. \quad (8)$$

The absolute values of expressions (7), (8) reflect the measure of the "dissimilarity" of the  $\mathbf{x}, \mathbf{y}$  intervals relative to their widths. In this case, the value of expression (8) is significantly less. This fact corresponds to the difference in the widths of the intervals  $\mathbf{x}, \mathbf{y}$  and makes it possible to build meaningful constructions on this basis, for example, in regularization procedures.

**Jaccard measure for sampling interval data.** The consistency measure introduced for two intervals in the form (3) admits a natural generalization to the case of interval sampling. Let there be an interval sample  $\mathbf{X} = \{\mathbf{x}_i\}, i = 1, 2, \dots, n$ . Let us define the measure  $JK(\mathbf{X})$  as

$$JK(\mathbf{X}) = \frac{\text{wid}(\Lambda_i \mathbf{x}_i)}{\text{wid}(\bigvee_i \mathbf{x}_i)}, \quad (9)$$

where

$$\bigvee_{1 \leq i \leq n} \mathbf{x}_i = \left[ \min_{1 \leq i \leq n} \underline{\mathbf{x}}_i, \max_{1 \leq i \leq n} \bar{\mathbf{x}}_i \right]; \quad \Lambda_{1 \leq i \leq n} \mathbf{x}_i = \left[ \max_{1 \leq i \leq n} \underline{\mathbf{x}}_i, \min_{1 \leq i \leq n} \bar{\mathbf{x}}_i \right].$$

It is important that in the case of an interval sample of two elements, expression (9) becomes expression (3).

Consistency measures for interval values are considered in publications [4–6], but only for consistent samples.

*Example.* Let us calculate the Jaccard coefficient (9) for consistent  $\mathbf{X}_1 = \{[0; 8], [2; 9], [3; 10]\}$  and inconsistent  $\mathbf{X}_2 = \{[9; 10], [10; 11], [11; 12]\}$  samples:

$$JK(\mathbf{X}_1) = \frac{8 - 3}{10 - 0} = \frac{5}{10} = 0.5,$$

$$JK(\mathbf{X}_2) = \frac{10 - 11}{12 - 9} = \frac{-1}{3} = -0.33(3).$$

**Jaccard measure and interval mode.** Let us consider the connection of the  $JK$  measure with other characteristics of interval samples, for example, with the interval mode of the sample. In the book [7], following the ideas of publication [8], this object (interval) is introduced as follows: "The interval sample mode  $\mathbf{X}$  is the set of intersection intervals of the largest consistent subsamples of the considered sample." The greatest length of consistent subsamples of a given sample is called the mode frequency.

Let us consider the simplest example illustrating the relationship between the Jaccard measure and the interval mode. In Fig. 2 we show options for calculating the interval mode and the mode consistency measure for a sample of two intervals  $\mathbf{X} = \{\mathbf{x}, \mathbf{y}\}$ , from which follows the relationship between the interval mode and the value of the Jaccard measure. In this case, the interval mode is a complex object that provides additional information about the sample. The interval mode is generally a multi-interval [7]. The nature of the Jaccard measure  $JK$  is different and has a numerical character. The numerator of expression (5) contains the inclusion minimum, and its position in multimode distributions may not belong to any sampling interval. In this case, additional expressions are required for the numerical description of the interval mode.

Preliminarily, it can be noted that the value of  $JK(\text{mode } \mathbf{X})$  in the case of a multi-interval mode is nonpositive. In this case, the absolute value of  $JK(\text{mode } \mathbf{X})$  depends both on the distance between the multisampling intervals and on the width of their inclusion maximum.

*Example.* For a noncovering sample  $\mathbf{X} = \{[1; 4], [5; 9], [1.5; 4.5], [6; 9]\}$ , the interval mode is the union of two intervals (multi-interval) mode  $\mathbf{X} = [1.5; 4] \cup [6; 9]$  and the measure of consistency is defined as

$$JK(\text{mode } \mathbf{X}) = \frac{4 - 6}{9 - 1} = -0.25.$$

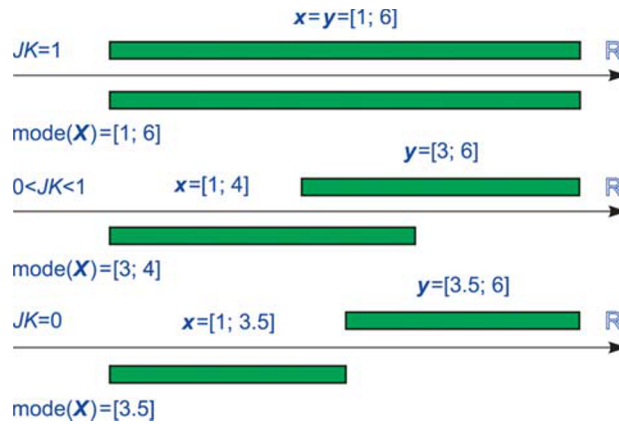


Fig. 2. Variants of relations between the position of intervals for the covering sample, interval mode and  $JK$  values.

In this case, the inclusion minimum  $\bigwedge_{1 \leq i \leq n} x_i = [6; 4]$  for the sample elements  $X$  in the numerator of expression (4) is an improper interval, which in the given case is also the inclusion minimum for the interval mode multi-interval elements.

**Application of  $JK$  to tween arithmetic and presentation of measurement data with errors of several kinds.** In practice, the ends of the intervals representing the results of measurements may not be known exactly, so that it becomes necessary to work with intervals that have interval ends. Such objects are known in interval analysis and are called tweens (English tween is an abbreviation for twice interval, double interval).

Tween, as an "interval of intervals" or an interval with interval ends, is represented as

$$X \quad [a, b] = \left[ [a, \bar{a}], [b, \bar{b}] \right]. \quad (10)$$

An alternative consideration of tweens as objects that give internal and external estimates of interval computations is proposed in [9]:

$$X = \left[ X^{\text{in}}, X^{\text{out}} \right] = \left[ \left[ \underline{X}^{\text{in}}, \bar{X}^{\text{in}} \right], \left[ \underline{X}^{\text{out}}, \bar{X}^{\text{out}} \right] \right], \quad (11)$$

where  $X^{\text{in}}$  and  $X^{\text{out}}$  are intervals of internal and external estimates of interval calculations.

In definitions (10), (11), a tween is defined by four values. The functional  $JK$  also depends on four values — the ends of two intervals  $x, y$ . If we take the internal and external components of the tween as arguments of  $JK$ , we obtain a relative value characterizing the relation of the internal and external estimates of the tween  $X$  in the form (11):

$$\delta X = JK \left( X^{\text{in}}, X^{\text{out}} \right). \quad (12)$$

*Example.* Given a temperature tween in the form (10):

$$T = \left[ [19.15; 19.85], [20.15; 20.85] \right] \text{ } ^\circ\text{C}.$$

We write it in the Nesterov form:

$$T = \left[ [19.85; 20.15], [19.15; 20.85] \right] \text{ } ^\circ\text{C}.$$

In accordance with construction (12), we have the "relative tween width"

$$\delta T = JK \left( T^{\text{in}}, T^{\text{out}} \right) = 0.18.$$

The positive value of  $\delta T$  reflects the consistency of the internal and external estimates of the interval object, and the absolute value is the "reliability margin" of the estimate. A negative value of  $\delta T$  corresponds to the nonempty intersection of the intervals of the lower and upper estimates of the interval  $T$ .

**Analysis of data with interval uncertainty.** Analysis of interval data emerged in the last decades of the 20th century as an alternative to traditional "probabilistic statistics" based on the methods of probability theory. Many works are devoted to various aspects of the analysis of interval data, in particular [10–19]. Analysis of data with interval uncertainty (statistics of interval data, analysis of interval data) is considered in the book [8], where the system of concepts and terms related to the processing of the specified data is presented. One of the most important concepts of interval statistics is a covering measurement (measurement) — an interval measurement result that is guaranteed to contain the true value of the measured value [7]. A measurement for which it is impossible to assert that it contains the true value of the measured quantity will be called noncovering. A covering measurement is a guaranteed two-sided "fork" of the values of the measured quantity, while for a noncovering measurement nothing of the kind can be asserted. This gives the covering measurements a fundamentally higher status and makes it possible to build more meaningful constructions on their basis.

Further, we will call a covering sample a set of measurements, the dominant part (most, etc.) of the measurements (observations) of which are covering. On the contrary, a sample is called noncovering if the predominant part of the measurements included in it are noncovering. Possible alternative terms are "inclusive dimension," "encompassing dimension" (their negation is "noninclusive", "nonencompassing"). Suggested English equivalents are enclosing measurement, covering measurement. Such a definition is not strict, which will avoid overly strict and rarely achievable sampling requirements.

When formulating new definitions, we follow the principle of correspondence in the methodology of science — any new scientific theory should include the old theory and its results as a special limiting case. As applied to the measure of consistency, further more general constructions contain (as special cases) the previously defined constructions. Symbolically, this can be expressed as

$$\mathbb{KR} \supseteq \mathbb{IR} \supseteq \mathbb{R} .$$

It is also important that the measures introduced for more general objects, namely, on samples of interval values and multidimensional interval objects, naturally reduce to constructions for simpler objects.

Thus, the measure  $JK(X)$  in the form (3) can be defined both for noncovering samples in Kaucher arithmetic and for covering samples in classical interval arithmetic, as well as for real "point" samples.

*Example.* We will demonstrate the use of the Jaccard index for data processing. The subject area refers to the physics of semiconductors — studies of the photoelectric characteristics of the sensor under test, carried out by specialists from the A. F. Ioffe Laboratory of Photoelectric Converters of the Physico-Technical Institute. There are two data samples with interval uncertainty.<sup>1</sup> Samples  $X_1, X_2$ , respectively, refer to sensors (photodetector, PD) PD1, PD2 (standard). The number of readings in the samples is 200. The PD1 sensor is calibrated according to the PD2 standard. The dependence between quantum efficiencies of the sensors  $QE_{PD1}, QE_{PD2}$  is assumed constant for each pair of measurement sets:

$$QE_{PD2} = \frac{I_{PD2}}{I_{PD1}} QE_{PD1} , \quad (13)$$

where  $I_{PD1}, I_{PD2}$  are the measured currents of the detectors.

Let us present the data in such a way as to apply the concepts of data statistics with interval uncertainty. One of the common ways to obtain interval results in primary measurements is the "intervaling" of point values, when an error interval  $\varepsilon$  is added to the point (base) value that is read from the indicators of the measuring device:

$$x = \dot{x} + \varepsilon .$$

We set the error interval as a balanced interval

$$\varepsilon = [-\varepsilon, \varepsilon] ,$$

<sup>1</sup>M. Z. Schwartz, working materials. URL: <https://github.com/AlexanderBazhenov/Solar-Data> (date of access: 11/18/2022).

for a specific procedure of measurements  $\varepsilon = 10^{-4}$  mV.

Due to a significant drift of the measuring equipment, the data is pre-processed and the drift is subtracted. After calculations, the data have a nonequivariant form, shown in Fig. 3.

According to the terminology of interval analysis, the sample under consideration is a vector of intervals, or an interval vector  $\mathbf{x} = (x_1, x_2, \dots, x_n)$ . In the case of estimating a single physical quantity from a sample of interval data, the information set will also be an information interval — an interval containing the values of the estimated quantity that are consistent with the measurements of the sample (consistent with the data of these measurements).

We calculate  $JK(\mathbf{X})$  for a set composed of the union of the values of the original samples  $\mathbf{X} = \{\mathbf{X}_1, \mathbf{X}_2\}$ . Using programs for calculations in interval arithmetic (cf. <https://github.com/szhilin/octave-interval-examples>), we obtain  $JK(\mathbf{X}) = -0.65$ . The negative value of  $JK(\mathbf{X})$  characterizes the inconsistency of the  $\mathbf{X}$  sample.

We use the model (13) of the connection between samples for different measurements ( $\mathbf{R}$  is an unknown factor):

$$\mathbf{X}_1 = \mathbf{R}\mathbf{X}_2.$$

In accordance with the tasks of researching sensors, it is necessary to give point and interval estimates of  $\mathbf{R}$ . We optimize the parameter  $\mathbf{R}$  as follows. A point estimate for the ratio of quantum efficiencies is

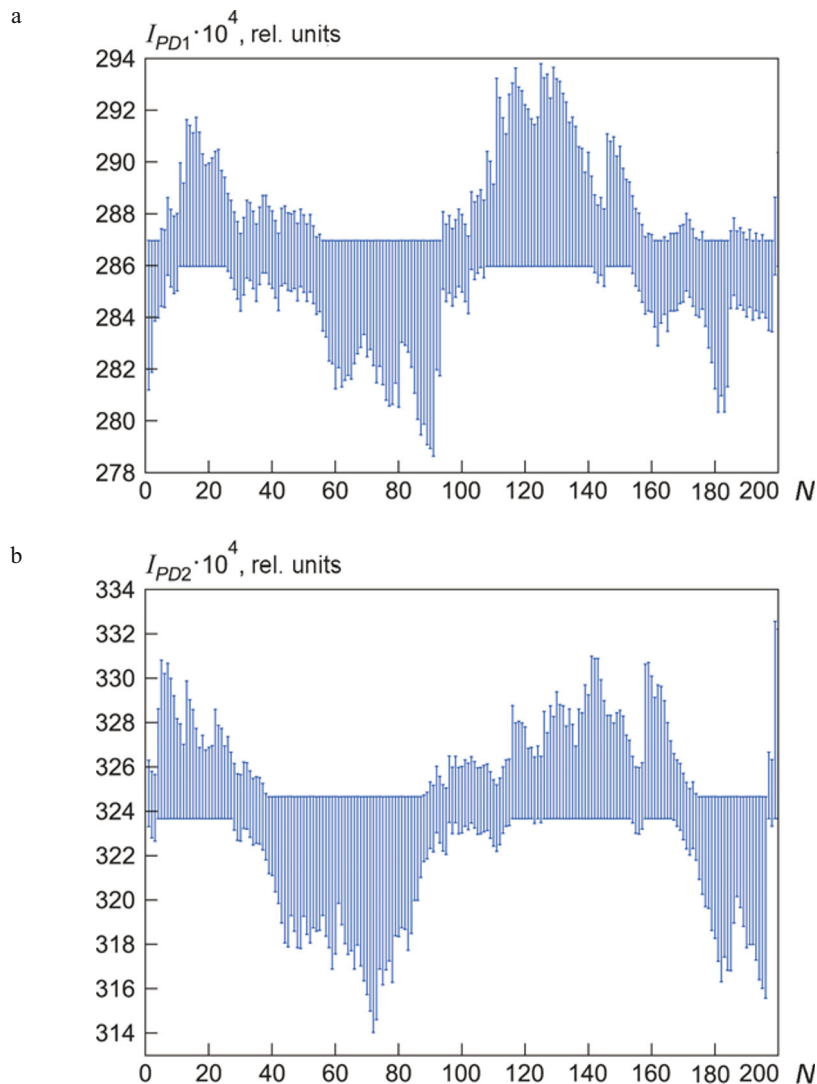


Fig. 3. Quantum efficiency measurement data samples (rel. units) for the studied (a) and reference (b) sensors.



$$\mathbf{R}_{\text{opt}} = \arg \max_{\mathbf{R}} JK(\mathbf{X}'), \quad (14)$$

where

$$\mathbf{X}' = \{\mathbf{R}\mathbf{X}_1, \mathbf{X}_2\}.$$

The interval estimate  $\mathbf{R}$  corresponds to the information set of the problem

$$\mathbf{R}: JK(\mathbf{X}') \geq 0. \quad (15)$$

Let us carry out calculations using formula (14) and plot  $JK(\mathbf{R})$  (Fig. 4), where the area corresponding to condition (15) is highlighted. The optimal value in formulation (14)  $\mathbf{R}_{\text{opt}} = 1.132$ , the information set of values (15)  $\mathbf{R} = [1.12; 1.14]$ .

The histogram of values of the  $\mathbf{X}'$  combined sample at  $\mathbf{R}_{\text{opt}}$  is shown in Fig. 5. The frequency distribution of histograms is close to symmetrical and unimodal, which positively indicates the correctness of the optimization procedure.

**Applications of measure of consistency.** For computational processes, two directions of application of the proposed measure can be proposed.

The first direction is the characterization of computational processes. When iteratively solving interval systems of equations, it is necessary to monitor the intersection of sets of successive operations. For this, a convenient means is an asymmetric measure equal to the ratio of the intersection of the sets of successive operations to the width of the iterative operator (Newton, Kravchik, etc.). The second direction is the solution of interval systems of linear algebraic equations (ISLAE). Traditionally, the solvability of ISLAE (in other words, the nonemptiness of solution sets) is studied using the technique of recognizing functionals. Recognizing functionals are constructed in such a way that their sign indicates that the solution sets are nonempty. For the main types of problem setting for solving ISLAE — finding a combined, admissible, and controlled decision sets — it is possible to build recognizing functionals based on symmetric and non-symmetric consistency measures.

For data analysis, the proposed measure can be used to characterize the range of a variable. When measuring any value with a set of sensors, it is necessary to allocate operating domains of parameters or measurement time ranges. A measure of consistency is naturally suitable for this. To select a reliable measurement range, it is necessary to set its practically reasonable threshold. A similar situation arises with multiple measurements of any quantity. For example, in the case of forward and reverse motion in the system under study, a hysteresis type of data may occur, including inconsistent ones. For characterization (identification of hysteresis), a measure of consistency can be used.

In clustering (factorization) problems, the working tool is the classification of data according to a set of features. A popular example is histogramming. If the data is essentially interval, one measurement may fall into more than one interval (bin) of the set according to which the classification is carried out. In this case, an asymmetric measure of consistency is well suited: the choice is made according to its maximum value.

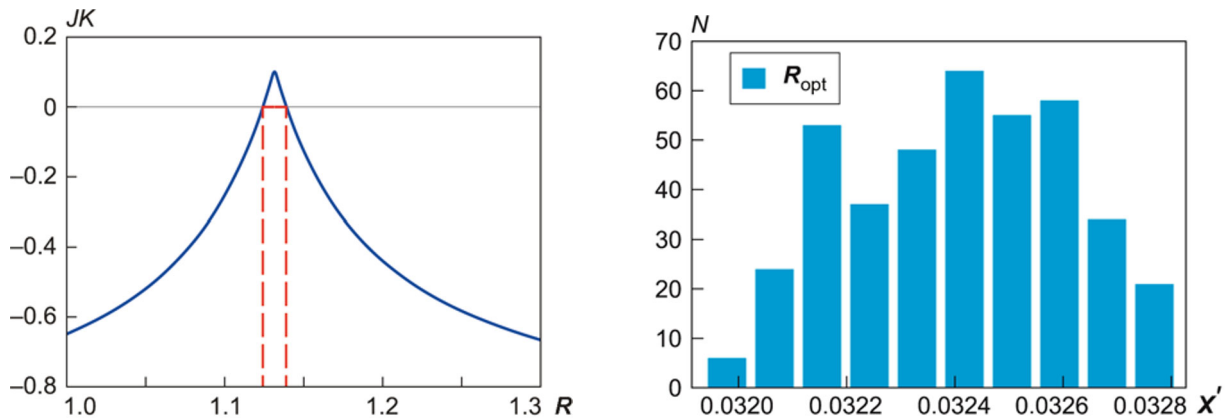


Fig. 4. Jaccard coefficient dependence on  $\mathbf{R}$ .

Fig. 5. Histogram of values of the  $\mathbf{X}'$  combined sample at  $\mathbf{R}_{\text{opt}}$ .



**Conclusion.** The proposed measure made it possible to generalize the concept of similarity of sets to interval objects: pairs of intervals and data samples with interval uncertainty. It is essential that the measure provides meaningful information for non-joint data samples as well. The above example demonstrated the possibility of obtaining an interval estimate of an unknown model parameter in the problem of dependency recovery. Directions for further research on the application of the proposed measure to characterize computational processes and analyze data with interval uncertainty are outlined.

**Acknowledgment.** The authors are grateful to the participants of the All-Russian webinar on interval analysis S. I. Zhilin, S. I. Kumkov, A. V. Prolubnikov, E. V. Chausova, and S. P. Shary for creative and constructive cooperation in the field of data analysis with interval uncertainty and also to V. M. Nesterov for discussion of issues related to twin arithmetic.

Works presented in the sections "Jaccard measure and interval mode", "Application  $JK$  to twin arithmetic and presentation of measurement data with several types of error", "Data analysis with interval uncertainty", supported by the Russian Science Foundation, project No. 21-72-20007. The works presented in the sections "Interval arithmetic", "Generalization of the Jaccard coefficient" and "Asymmetric variants of the  $JK(x, y)$  measure" were supported by the Ioffe Institute RAS (within the framework of the RF state task 0034-2019-0001 and 0040-2019-0023).

## REFERENCES

1. B. I. Semkin, On the Relation Between Mean Values of Two Measures of Inclusion and Measures of Similarity, *Byull. Botanicheskogo sada-institutu DVO RAS*, No. 3, 91–101 (2009).
2. R. B. Kearfott, M. T. Nakao, A. Neumaier, S. M. Rump, S. P. Shary, and P. van Hentenryck, Standardized Notation in Interval Analysis, *Comput. Technol.*, **15**, No. 1, 7–13 (2010).
3. S. Shary, *Numerical Computation of Formal Solutions to Interval Linear Systems of Equations*, arXiv:1903.10272v1 [math.NA], <https://doi.org/10.48550/arXiv.1903.10272>.
4. S. Kabir, C. Wagner, T. C. Havens, D. T. Anderson, and U. Aickelin, *IEEE Int. Conf. on Fuzzy Systems (FUZZ-IEEE 2017)*, 2017, <https://doi.org/10.1109/FUZZ-IEEE.2017.8015623>.
5. T. Wilkin and G. Beliakov, *IEEE Int. Conf. on Fuzzy Systems (FUZZ-IEEE 2019)*, 1–6 (2019), <https://doi.org/10.1109/FUZZ-IEEE.2019.8858850>.
6. S. Kabir, C. Wagner, and Z. Ellerby, *Towards Handling Uncertainty-at-Source in AI — A Review and Next Steps for Interval Regression*, arXiv:2104.07245 [cs.LG], <https://doi.org/10.48550/arXiv.2104.07245>.
7. A. N. Bazhenov, S. I. Zhilin, S. I. Kumkov, and S. P. Sharyj, *Processing and Analysis of Data with Interval Uncertainty*, 2022, available at: <http://www.nsc.ru/interval/Library/AppBooks/InteData Processing.pdf> (accessed: 10.11.2022).
8. C. Hu and Z. H. Hu, On Statistics, Probability, and Entropy of Interval-Valued Datasets, *Information Processing and Management of Uncertainty in Knowledge-Based Systems. IPMU 2020. Communications in Computer and Information Science*, M. J. Lesot et al. Eds., Cham, Springer, **1239**, 2020, [https://doi.org/10.1007/978-3-030-50153-2\\_31](https://doi.org/10.1007/978-3-030-50153-2_31).
9. V. M. Nesterov, *Tween Arithmetics and Their Application in Methods and Algorithms of Two-Sided Interval Estimation*, *Diss. Doct. Phis. Math. Sci. St. Petersburg (St. Petersburg Institute of Informatics and Automation of the Russian Academy of Sciences)*, 1999.
10. S. Shary, *Comput. Technol.*, **2**, No. 2, 150–172 (2017), <http://dx.doi.org/10.14529/mmph170105>.
11. S. Shary, *J. Comput. Syst. Sci. Int.*, **56**, No. 6, 897–913 (2017), <https://doi.org/10.7868/S0002338817060014>.
12. S. Shary, Identification of Outliers in the Maximum Mat-Ching Method in the Analysis of Interval Data, *Proc. All-Russian Conf. on Mathematics Int. Participation "MAC-2018"*, Barnaul, AltGU Publishing House, 2018, pp. 215–218.
13. S. Shary, On a Variability Measure for Estimates of Parameters in the Statistics of Interval Data, *Comput. Technol.*, **24**, No. 5, 90–108 (2019), <https://doi.org/10.25743/ICT.2019.24.5.008>.
14. S. P. Shary, Data Fitting Problem under Interval Uncertainty in Data, *Ind. Lab. Diagn. Mater.*, **86**, No. 1, 62–74, (2020), <https://doi.org/10.26896/1028-6861-2020-86-1-62-74>.
15. S. I. Zhilin, *Reliab. Comput.*, **11**, 433–442 (2005), <https://doi.org/10.1007/s11155-005-0050-3>.
16. S. I. Zhilin, *Chemometr. Intell. Lab. Syst.*, **88**, No. 1, 60–68 (2007), <https://doi.org/10.1016/j.chemolab.2006.10.004>.
17. S. I. Kumkov, Processing of Experimental Data on Ionic Conductivity of Molten Electrolyte by Methods of Interval Analysis, *Russian Metallurgy (METALLY)*, No. 3, 79–89 (2010).
18. S. I. Kumkov and Yu. V. Mikushina, *Reliab. Comput.*, **19**, 197–214 (2013).
19. H. T. Nguyen, V. Kreinovich, B. Wu, and G. Xiang, *Computing Statistics under Interval and Fuzzy Uncertainty. Applications to Computer Science and Engineering*, Springer, Berlin-Heidelberg, 2012, <https://doi.org/10.1007/978-3-642-24905-1>.