
Interval-Valued and Fuzzy-Valued Random Variables: From Computing Sample Variances to Computing Sample Covariances

Jan B. Beck¹, Vladik Kreinovich¹, and Berlin Wu²

¹ University of Texas, El Paso TX 79968, USA, {janb,vladik}@cs.utep.edu

² National Chengchi University, Taipei, Taiwan berlin@math.nccu.edu.tw

Summary. Due to measurement uncertainty, often, instead of the actual values x_i of the measured quantities, we only know the intervals $\mathbf{x}_i = [\tilde{x}_i - \Delta_i, \tilde{x}_i + \Delta_i]$, where \tilde{x}_i is the measured value and Δ_i is the upper bound on the measurement error (provided, e.g., by the manufacturer of the measuring instrument). These intervals can be viewed as *random intervals*, i.e., as samples from the interval-valued random variable. In such situations, instead of the exact value of the sample statistics such as covariance $C_{x,y}$, we can only have an interval $\mathbf{C}_{x,y}$ of possible values of this statistic. It is known that in general, computing such an interval $\mathbf{C}_{x,y}$ for $C_{x,y}$ is an NP-hard problem. In this paper, we describe an algorithm that computes this range $\mathbf{C}_{x,y}$ for the case when the measurements are accurate enough – so that the intervals corresponding to different measurements do not intersect much.

Introduction

Traditional statistics: brief reminder. In traditional statistics, we have i.i.d. variables x_1, \dots, x_n, \dots , and we want to find statistics $s_n(x_1, \dots, x_n)$ that would approximate the desired parameter of the corresponding probability distribution. For example, if we want to estimate the mean $E[x]$, we can take the arithmetic average $s_n = (x_1 + \dots + x_n)/n$. It is known that as $n \rightarrow \infty$, this statistic tends (with probability 1) to the desired mean: $s_n \rightarrow E[x]$. Similarly, the sample variance tends to the actual variance, the sample covariance between two different samples tends to the actual covariance $C[x, y]$, etc.

Coarsening: a source of random sets. In traditional statistics, we implicitly assume that the values x_i are directly observable. In real life, due to (inevitable) measurement uncertainty, often, what we actually observe is a *set* S_i that contains the actual (unknown) value of x_i . This phenomenon is called *coarsening*; see, e.g., [3]. Due to coarsening, instead of the actual values x_i , all we know is the sets X_1, \dots, X_n, \dots that are known to contain the actual (un-observable) values x_i : $x_i \in X_i$.

Statistics based on coarsening. The sets X_1, \dots, X_n, \dots are i.i.d. *random sets*. We want to find statistics of these random sets that would enable us to approximate the desired parameters of the original distribution x . Here, a statistic $S_n(X_1, \dots, X_n)$ transform n sets X_1, \dots, X_n into a new set. We want this statistic $S_n(X_1, \dots, X_n)$ to tend to a limit set L as $n \rightarrow \infty$, and we want this limit set L to contain the value of the desired parameter of the original distribution.

For example, if we are interested in the mean $E[x]$, then we can take $S_n = (X_1 + \dots + X_n)/n$ (where the sum is the Mankowski – element-wise – sum of the sets). It is possible to show that, under reasonable assumptions, this statistic tends to a limit L , and that $E[x] \in L$. This limit can be viewed, therefore, as a set-based average of the sets X_1, \dots, X_n .

Important issue: computational complexity. There has been a lot of interesting theoretical research on set-valued random variables and corresponding statistics. In many cases, the corresponding statistics have been designed, and their asymptotical properties have been proven; see, e.g., [2, 6] and references therein.

In many such situations, the main obstacle on the way of practically using these statistics is the fact that going from random numbers to random sets drastically increases the computational complexity – hence, the running time – of the required computations. It is therefore desirable to come up with new, faster algorithms for computing such set-values heuristics.

What we are planning to do. In this paper, we design such faster algorithm for a specific practical problem: the problem of data processing. As we will show, in this problem, it is reasonable to restrict ourselves to *random intervals* – an important specific case of random sets.

Practical Problem: Data Processing – From Computing to Probabilities to Intervals

Why data processing? In many real-life situations, we are interested in the value of a physical quantity y that is difficult or impossible to measure directly. Examples of such quantities are the distance to a star and the amount of oil in a given well. Since we cannot measure y directly, a natural idea is to measure y *indirectly*. Specifically, we find some easier-to-measure quantities x_1, \dots, x_n which are related to y by a known relation $y = f(x_1, \dots, x_n)$. This relation may be a simple functional transformation, or complex algorithm (e.g., for the amount of oil, numerical solution to an inverse problem).

It is worth mentioning that in the vast majority of these cases, the function $f(x_1, \dots, x_n)$ that describes the dependence between physical quantities is continuous.

Then, to estimate y , we first measure the values of the quantities x_1, \dots, x_n , and then we use the results $\tilde{x}_1, \dots, \tilde{x}_n$ of these measurements to compute an estimate \tilde{y} for y as $\tilde{y} = f(\tilde{x}_1, \dots, \tilde{x}_n)$.

For example, to find the resistance R , we measure current I and voltage V , and then use the known relation $R = V/I$ to estimate resistance as $\tilde{R} = V/I$.

Computing an estimate for y based on the results of direct measurements is called *data processing*; data processing is the main reason why computers were invented in the first place, and data processing is still one of the main uses of computers as number crunching devices.

Comment. In this paper, for simplicity, we consider the case when the relation between x_i and y is known exactly; in some practical situations, we only know an approximate relation between x_i and y .

Why interval computations? From computing to probabilities to intervals. Measurement are never 100% accurate, so in reality, the actual value x_i of i -th measured quantity can differ from the measurement result \tilde{x}_i . Because of these *measurement errors* $\Delta x_i \stackrel{\text{def}}{=} \tilde{x}_i - x_i$, the result $\tilde{y} = f(\tilde{x}_1, \dots, \tilde{x}_n)$ of data processing is, in general, different from the actual value $y = f(x_1, \dots, x_n)$ of the desired quantity y [12].

It is desirable to describe the error $\Delta y \stackrel{\text{def}}{=} \tilde{y} - y$ of the result of data processing. To do that, we must have some information about the errors of direct measurements.

What do we know about the errors Δx_i of direct measurements? First, the manufacturer of the measuring instrument must supply us with an upper bound Δ_i on the measurement error. If no such upper bound is supplied, this means that no accuracy is guaranteed, and the corresponding “measuring instrument” is practically useless. In this case, once we performed a measurement and got a measurement result \tilde{x}_i , we know that the actual (unknown) value x_i of the measured quantity belongs to the interval $\mathbf{x}_i = [\underline{x}_i, \bar{x}_i]$, where $\underline{x}_i = \tilde{x}_i - \Delta_i$ and $\bar{x}_i = \tilde{x}_i + \Delta_i$.

In many practical situations, we not only know the interval $[-\Delta_i, \Delta_i]$ of possible values of the measurement error; we also know the probability of different values Δx_i within this interval. This knowledge underlies the traditional engineering approach to estimating the error of indirect measurement, in which we assume that we know the probability distributions for measurement errors Δx_i .

In practice, we can determine the desired probabilities of different values of Δx_i by comparing the results of measuring with this instrument with the results of measuring the same quantity by a standard (much more accurate) measuring instrument. Since the standard measuring instrument is much more accurate than the one use, the difference between these two measurement results is practically equal to the measurement error; thus, the empirical distribution of this difference is close to the desired probability distribution for measurement error. There are two cases, however, when this determination is not done:

- First is the case of cutting-edge measurements, e.g., measurements in fundamental science. When a Hubble telescope detects the light from a distant

galaxy, there is no “standard” (much more accurate) telescope floating nearby that we can use to calibrate the Hubble: the Hubble telescope is the best we have.

- The second case is the case of measurements on the shop floor. In this case, in principle, every sensor can be thoroughly calibrated, but sensor calibration is so costly – usually costing ten times more than the sensor itself – that manufacturers rarely do it.

In both cases, we have no information about the probabilities of Δx_i ; the only information we have is the upper bound on the measurement error.

In this case, after we performed a measurement and got a measurement result \tilde{x}_i , the only information that we have about the actual value x_i of the measured quantity is that it belongs to the interval $\mathbf{x}_i = [\tilde{x}_i - \Delta_i, \tilde{x}_i + \Delta_i]$. In such situations, the only information that we have about the (unknown) actual value of $y = f(x_1, \dots, x_n)$ is that y belongs to the range $\mathbf{y} = [\underline{y}, \bar{y}]$ of the function f over the box $\mathbf{x}_1 \times \dots \times \mathbf{x}_n$:

$$\mathbf{y} = [\underline{y}, \bar{y}] = \{f(x_1, \dots, x_n) \mid x_1 \in \mathbf{x}_1, \dots, x_n \in \mathbf{x}_n\}.$$

For continuous functions $f(x_1, \dots, x_n)$, this range is an interval. The process of computing this interval range based on the input intervals \mathbf{x}_i is called *interval computations*; see, e.g., [4].

Interval computations techniques: brief reminder. Historically the first method for computing the enclosure for the range is the method which is sometimes called “straightforward” interval computations. This method is based on the fact that inside the computer, every algorithm consists of elementary operations (arithmetic operations, min, max, etc.). For each elementary operation $f(a, b)$, if we know the intervals \mathbf{a} and \mathbf{b} for a and b , we can compute the exact range $f(\mathbf{a}, \mathbf{b})$. The corresponding formulas form the so-called *interval arithmetic*. For example,

$$[\underline{a}, \bar{a}] + [\underline{b}, \bar{b}] = [\underline{a} + \underline{b}, \bar{a} + \bar{b}]; \quad [\underline{a}, \bar{a}] - [\underline{b}, \bar{b}] = [\underline{a} - \bar{b}, \bar{a} - \underline{b}];$$

$$[\underline{a}, \bar{a}] \cdot [\underline{b}, \bar{b}] = [\min(\underline{a} \cdot \underline{b}, \underline{a} \cdot \bar{b}, \bar{a} \cdot \underline{b}, \bar{a} \cdot \bar{b}), \max(\underline{a} \cdot \underline{b}, \underline{a} \cdot \bar{b}, \bar{a} \cdot \underline{b}, \bar{a} \cdot \bar{b})].$$

In straightforward interval computations, we repeat the computations forming the program f step-by-step, replacing each operation with real numbers by the corresponding operation of interval arithmetic. It is known that, as a result, we get an enclosure $\mathbf{Y} \supseteq \mathbf{y}$ for the desired range.

In some cases, this enclosure is exact. In more complex cases (see examples below), the enclosure has excess width.

There exist more sophisticated techniques for producing a narrower enclosure, e.g., a centered form method. However, for each of these techniques, there are cases when we get an excess width. Reason: as shown in [5], the problem of computing the exact range is known to be NP-hard even for polynomial functions $f(x_1, \dots, x_n)$ (actually, even for quadratic functions f).

Comment. NP-hard means, crudely speaking, that no feasible algorithm can compute the exact range of $f(x_1, \dots, x_n)$ for all possible polynomials $f(x_1, \dots, x_n)$ and for all possible intervals $\mathbf{x}_1, \dots, \mathbf{x}_n$.

What we are planning to do? In this paper, we analyze a specific interval computations problem – when we use traditional statistical data processing algorithms $f(x_1, \dots, x_n)$ to process the results of direct measurements.

Error Estimation for Traditional Statistical Data Processing Algorithms under Interval Uncertainty: Known Results

Formulation of the problem. When we have n results x_1, \dots, x_n of repeated measurement of the same quantity (at different points, or at different moments of time), traditional statistical approach usually starts with computing their sample average $E_x = (x_1 + \dots + x_n)/n$ and their (sample) variance

$$V_x = \frac{(x_1 - E_x)^2 + \dots + (x_n - E_x)^2}{n} \quad (1)$$

(or, equivalently, the sample standard deviation $\sigma = \sqrt{V}$). If, during each measurement i , we measure the values x_i and y_i of two different quantities x and y , then we also compute their (sample) covariance

$$C_{x,y} = \frac{(x_1 - E_x) \cdot (y_1 - E_y) + \dots + (x_n - E_x) \cdot (y_n - E_y)}{n}, \quad (2)$$

see, e.g., [12].

As we have mentioned, in real life, we often do not know the exact values of the quantities x_i and y_i , we only know the intervals \mathbf{x}_i of possible values of x_i and the intervals \mathbf{y}_i of possible values of y_i . In such situations, for different possible values $x_i \in \mathbf{x}_i$ and $y_i \in \mathbf{y}_i$, we get different values of E_x , E_y , V_x , and $C_{x,y}$. The question is: what are the intervals \mathbf{E}_x , \mathbf{V}_x , and $\mathbf{C}_{x,y}$ of possible values of E_x , V_x , and $C_{x,y}$?

The practical importance of this question was emphasized, e.g., in [9, 10] on the example of processing geophysical data.

Comment: the problem reformulated in terms of set-valued random variables. Traditional statistical data processing means that we assume that the measured values x_i and y_i are samples of random variables, and based on these samples, we are estimating the actual average, variance, and covariance.

Similarly, in case of interval uncertainty, we can say that the intervals \mathbf{x}_i and \mathbf{y}_i coming from measurements are samples of interval-valued random variables, and we are interested in estimating the actual (properly defined) average, variance, and covariance of these interval-valued random variables.

Bounds on E. For E_x , the straightforward interval computations leads to the exact range:

$$\mathbf{E}_x = \frac{\mathbf{x}_1 + \dots + \mathbf{x}_n}{n}, \text{ i.e., } \underline{E}_x = \frac{\underline{x}_1 + \dots + \underline{x}_n}{n}, \text{ and } \overline{E}_x = \frac{\overline{x}_1 + \dots + \overline{x}_n}{n}.$$

For variance, the problem is difficult. For \mathbf{V}_x , all known algorithms lead to an excess width. Specifically, there exist feasible algorithms for computing \underline{V}_x (see, e.g., [1]), but in general, the problem of computing \overline{V}_x is NP-hard [1].

It is also known that in some practically important cases, feasible algorithms for computing \overline{V}_x are possible. One such practically useful case is when the measurement accuracy is good enough so that we can tell that the different measured values \tilde{x}_i are indeed different – e.g., the corresponding intervals \mathbf{x}_i do not intersect. In this case, there exists a quadratic-time algorithm for computing \overline{V}_x ; see, e.g., [1].

What about covariance? The only thing that we know is that in general, computing covariance $\mathbf{C}_{x,y}$ is NP-hard [11].

In this paper, we show that (similarly to the case of variance), it is possible to compute the interval covariance when the measurement are accurate enough to enable us to distinguish between different measurement results $(\tilde{x}_i, \tilde{y}_i)$.

Main Results

Theorem 1. *There exists a polynomial-time algorithm that, given a list of n pairwise disjoint boxes $\mathbf{x}_i \times \mathbf{y}_i$ ($1 \leq i \leq n$) (i.e., in which every two boxes have an empty intersection), produces the exact range $\mathbf{C}_{x,y}$ for the covariance $C_{x,y}$.*

Theorem 2. *For every integer $K > 1$, there exists a polynomial-time algorithm that, given a list of n boxes $\mathbf{x}_i \times \mathbf{y}_i$ ($1 \leq i \leq n$) in which $> K$ boxes always have an empty intersection, produces the exact range $\mathbf{C}_{x,y}$ for the covariance $C_{x,y}$.*

Proof. Since $C_{x,y}$ is linear in x_i , we have $C_{-x,y} = -C_{x,y}$ hence $\mathbf{C}_{-x,y} = -\mathbf{C}_{x,y}$, so $\underline{C}_{-x,y} = -\overline{C}_{x,y}$. Because of this relation, it is sufficient to provide an algorithm for computing $\underline{C}_{x,y}$; we will then compute $\overline{C}_{x,y}$ as $-\underline{C}_{-x,y}$.

The function $C_{x,y}$ is linear in each of its variables x_i and y_i . In general, a linear function $f(x)$ attains its minimum on an interval $[\underline{x}, \overline{x}]$ at one of its endpoints: at \underline{x} if f is non-decreasing ($\partial f / \partial x \geq 0$) and at \overline{x} if f is non-increasing ($\partial f / \partial x \leq 0$). For $C_{x,y}$, we have $\partial C_{x,y} / \partial x_i = (1/n) \cdot y_i - (1/n) \cdot E_y$, so $\partial C_{x,y} / \partial x_i \geq 0$ if and only if $y_i \geq E_y$. Thus, for each i , the values x_i^m and y_i^m at which $C_{x,y}$ attains its minimum satisfy the following four properties:

1. if $x_i^m = \underline{x}_i$, then $y_i^m \geq E_y$;
2. if $x_i^m = \overline{x}_i$, then $y_i^m \leq E_y$;
3. if $y_i^m = \underline{y}_i$, then $x_i^m \geq E_x$;
4. if $y_i^m = \overline{y}_i$, then $x_i^m \leq E_x$.

Let us show that if we know the vector $E \stackrel{\text{def}}{=} (E_x, E_y)$, and this vector is outside the i -th box $\mathbf{b}_i \stackrel{\text{def}}{=} \mathbf{x}_i \times \mathbf{y}_i$, then we can uniquely determine the values x_i^m and y_i^m .

Indeed, the fact that $E \notin \mathbf{b}_i$ means that either $E_x \notin \mathbf{x}_i$ or $E_y \notin \mathbf{y}_i$. Without losing generality, let us assume that $E_x \notin \mathbf{x}_i$, i.e., that either $E_x < \underline{x}_i$ or $E_x > \overline{x}_i$.

If $E_x < \underline{x}_i$, then, since $\underline{x}_i \leq x_i^m$, we have $E_x < x_i^m$. Hence, according to Property 4, we cannot have $y_i^m = \bar{y}_i$. Since the minimum is always attained at one of the endpoints, we thus have $y_i^m = \underline{y}_i$. Now that we know the value of y_i^m , we can use Properties 1 and 2:

$$\text{if } \underline{y}_i \geq E_y, \text{ then } x_i^m = \underline{x}_i; \quad \text{if } \underline{y}_i \leq E_y, \text{ then } x_i^m = \underline{x}_i.$$

Similarly, if $E_x > \bar{x}_i$, then $y_i^m = \bar{y}_i$, and:

$$\text{if } \bar{y}_i \geq E_y, \text{ then } x_i^m = \underline{x}_i; \quad \text{if } \bar{y}_i \leq E_y, \text{ then } x_i^m = \underline{x}_i.$$

So, to compute $\underline{C}_{x,y}$, we sort all $2n$ values $\underline{x}_i, \bar{x}_i$ into a sequence $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(2n)}$, and we sort the $2n$ values $\underline{y}_i, \bar{y}_i$ into a sequence $y_{(1)} \leq y_{(2)} \leq \dots \leq y_{(2n)}$. We thus get $2n \times 2n$ "zones" $\mathbf{z}_{k,l} \stackrel{\text{def}}{=} [x_{(k)}, x_{(k+1)}] \times [y_{(l)}, x_{(l+1)}]$.

We know that the average E of the actual minimum values is attained in one of these zones. If we assume that $E \in \mathbf{z}_{k,l}$, i.e., in particular, that $E_x \in [x_{(k)}, x_{(k+1)}]$, then the condition $\underline{x}_i \geq E_x$ is guaranteed to be satisfied if $\underline{x}_i \geq x_{(k+1)}$. Thus, following the above arguments, we can find the values (x_i^m, y_i^m) for all the boxes \mathbf{b}_i that do not contain this zone:

	$\bar{y}_i \leq y_{(l)}$	$\underline{y}_i \leq y_{(l)} \leq y_{(l+1)} \leq \bar{y}_i$	$y_{(l+1)} \leq \underline{y}_i$
$\bar{x}_i \leq x_{(k)}$	(\bar{x}_i, \bar{y}_i)	$(\underline{x}_i, \bar{y}_i)$	$(\underline{x}_i, \bar{y}_i)$
$\underline{x}_i \leq x_{(k)} \leq x_{(k+1)} \leq \bar{x}_i$	$(\bar{x}_i, \underline{y}_i)$?	$(\underline{x}_i, \bar{y}_i)$
$x_{(k+1)} \leq \underline{x}_i$	$(\bar{x}_i, \underline{y}_i)$	$(\bar{x}_i, \underline{y}_i)$	$(\underline{x}_i, \underline{y}_i)$

As we can see, for each of $O(n^2)$ zones $\mathbf{z}_{k,l}$, the only case when we do not know the corresponding values (x_i^m, y_i^m) is when \mathbf{b}_i contains this zone. All boxes \mathbf{b}_i with this property have a common intersection $\mathbf{z}_{k,l}$, thus, there can be no more than K of them. For each of these $\leq K$ boxes \mathbf{b}_i , we try all 4 possible combinations of endpoints as the corresponding (x_i^m, y_i^m) .

Thus, for each of $O(n^2)$ zones, we must try $\leq 4^K$ possible sequences of pairs (x_i^m, y_i^m) . We compute each of these n -element sequences element-by-element, so computing each sequence requires $O(n)$ computational steps.

For each of these sequences, we check whether the averages E_x and E_y are indeed within this zone, and if they are, we compute the correlation. The smallest of the resulting correlations is the desired value $\underline{C}_{x,y}$.

For each of $O(n^2)$ zones, we need $O(n)$ steps, to the total of $O(n^2) \times O(n) = O(n^3)$; computing the smallest of $O(n^2)$ values requires $O(n^2)$ more steps. Thus, our algorithm computes $\underline{C}_{x,y}$ in $O(n^3)$ steps. \square

From Interval-Valued to Fuzzy-Valued Random Variables

Often, in addition (or instead) the guaranteed bounds, an expert can provide bounds that contain x_i with a certain degree of confidence. Often, we know several such bounding intervals corresponding to different degrees of confidence. Such a nested family of intervals is also called a *fuzzy set*, because it

turns out to be equivalent to a more traditional definition of fuzzy set [7, 8] (if a traditional fuzzy set is given, then different intervals from the nested family can be viewed as α -cuts corresponding to different levels of uncertainty α).

To provide statistical values of fuzzy-valued random variables, we can therefore, for each level α , apply the above interval-valued techniques to the corresponding α -cuts.

Acknowledgments. This work was supported by NASA grant NCC5-209, by the Air Force Office of Scientific Research grant F49620-00-1-0365, by NSF grants EAR-0112968, EAR-0225670, and EIA-0321328, and by the Army Research Laboratories grant DATM-05-02-C-0046. The authors are thankful to Hung T. Nguyen and to the anonymous referees for valuable suggestions.

References

1. Ferson S, Ginzburg L, Kreinovich V, Longpré L, Aviles M (2002) Computing Variance for Interval Data is NP-Hard, ACM SIGACT News 33(2):108–118.
2. Goutsias J, Mahler RPS, Nguyen HT, eds. (1997) Random Sets: Theory and Applications. Springer-Verlag, N.Y.
3. Heitjan DF, Rubin DB (1991) Ignorability and coarse data, Ann. Stat. 19(4):2244–2253
4. Jaulin L, Kieffer M, Didrit O, and Walter E (2001) Applied Interval Analysis. Springer-Verlag, Berlin
5. Kreinovich V, Lakeyev A, Rohn J, Kahl P (1997) Computational Complexity and Feasibility of Data Processing and Interval Computations. Kluwer, Dordrecht
6. Li S, Ogura Y, Kreinovich V (2002) Limit Theorems and Applications of Set Valued and Fuzzy Valued Random Variables. Kluwer, Dordrecht
7. Nguyen HT, Kreinovich V (1996) Nested Intervals and Sets: Concepts, Relations to Fuzzy Sets, and Applications, In: Kearfott RB et al., Applications of Interval Computations, Kluwer, Dordrecht, 245–290
8. Nguyen HT, Walker EA (1999) First Course in Fuzzy Logic, CRC Press, Boca Raton, Florida
9. Nivlet P, Fournier F, and Royer J. (2001) A new methodology to account for uncertainties in 4-D seismic interpretation, Proc. 71st Annual Int'l Meeting of Soc. of Exploratory Geophysics SEG'2001, San Antonio, TX, September 9–14, 1644–1647.
10. Nivlet P, Fournier F, Royer J (2001) Propagating interval uncertainties in supervised pattern recognition for reservoir characterization, Proc. 2001 Society of Petroleum Engineers Annual Conf. SPE'2001, New Orleans, LA, September 30–October 3, paper SPE-71327.
11. Osegueda R, Kreinovich V, Potluri L, Aló R (2002) Non-Destructive Testing of Aerospace Structures: Granularity and Data Mining Approach. In: Proc. FUZZ-IEEE'2002, Honolulu, HI, May 12–17, 2002, Vol. 1, 685–689
12. Rabinovich S (1993) Measurement Errors: Theory and Practice. American Institute of Physics, New York