

Exact Bounds on Sample Variance of Interval Data

Scott Ferson¹, Lev Ginzburg¹,
Vladik Kreinovich², Luc Longpré, and Monica Aviles²

¹Applied Biomathematics, 100 North Country Road,
Setauket, NY 11733, USA, {scott,lev}@ramas.com

²Computer Science Department, University of Texas at El Paso
El Paso, TX 79968, USA, {maviles,longpre,vladik}@cs.utep.edu

Abstract

We provide a feasible (quadratic time) algorithm for computing the lower bound \underline{V} on the sample variance of interval data. We prove that the problem of computing the upper bound \overline{V} is, in general, NP-hard. We provide a feasible algorithm that computes \overline{V} under reasonable easily verifiable conditions.

1 Introduction

1.1 Formulation of the Problem

When we have n measurement results x_1, \dots, x_n , traditional statistical approach usually starts with computing their sample average

$$E = \bar{x} = \frac{x_1 + \dots + x_n}{n}$$

and their sample variance

$$V = \frac{(x_1 - E)^2 + \dots + (x_n - E)^2}{n - 1} \quad (1.1)$$

(or, equivalently, the sample standard deviation $\sigma = \sqrt{V}$); see, e.g., [9].

Sample variance is an unbiased estimator of the variance of the distribution from which observations are assumed to be randomly sampled. For Gaussian distribution, this estimator is a maximum likelihood estimator of the distribution variance.

In some practical situations, we only have intervals $\mathbf{x}_i = [\underline{x}_i, \overline{x}_i]$ of possible values of x_i . This happens, for example, if instead of observing the actual value x_i of the random variable, we observe the value \tilde{x}_i measured by an instrument with a known upper bound Δ_i on the measurement error; then, the actual (unknown) value is within the interval $\mathbf{x}_i = [\tilde{x}_i - \Delta_i, \tilde{x}_i + \Delta_i]$.

As a result, the sets of possible values of E and V are also intervals. The interval \mathbf{E} for the sample average can be obtained by using straightforward interval computations, i.e., by replacing each elementary operation with numbers by the corresponding operation of interval arithmetic:

$$\mathbf{E} = \frac{\mathbf{x}_1 + \dots + \mathbf{x}_n}{n}. \quad (1.2)$$

What is the interval $[\underline{V}, \overline{V}]$ of possible values for sample variance V ?

When all the intervals \mathbf{x}_i intersect, then it is possible that all the actual (unknown) values $x_i \in \mathbf{x}_i$ are the same and hence, that the sample variance is 0. In other words, if the intervals have a non-empty intersection, then $\underline{V} = 0$. Conversely, if the intersection of \mathbf{x}_i is empty, then V cannot be 0, hence $\underline{V} > 0$. The question is (see, e.g., [13]): What is the total set of possible values of V when the above intersection is empty?

1.2 For this Problem, Traditional Interval Methods Sometimes Overestimate

Let us show that for this problem, traditional interval methods sometimes overestimate.

1.2.1 Straightforward Interval Computations

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 $f(x, y)$, if we know the intervals \mathbf{x} and \mathbf{y} for x and y , we can compute the exact range $f(\mathbf{x}, \mathbf{y})$. The corresponding formulas form the so-called *interval arithmetic*. 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 for the desired range.

For the problem of computing the range of sample variance, straightforward interval computations sometimes overestimate. For example, for $\mathbf{x}_1 = \mathbf{x}_2 = [0, 1]$, the actual $V = (x_1 - x_2)^2/2$ and hence, the actual range $\mathbf{V} = [0, 0.5]$. On the other hand, $\mathbf{E} = [0, 1]$, hence

$$(\mathbf{x}_1 - \mathbf{E})^2 + (\mathbf{x}_2 - \mathbf{E})^2 = [0, 2] \supset [0, 0.5].$$

1.2.2 Centered Form

A better range is often provided by a *centered form*, in which a range $f(\mathbf{x}_1, \dots, \mathbf{x}_n)$ of a smooth function on a box $\mathbf{x}_1 \times \dots \times \mathbf{x}_n$ is estimated as

$$f(\mathbf{x}_1, \dots, \mathbf{x}_n) \subseteq f(\tilde{x}_1, \dots, \tilde{x}_n) + \sum_{i=1}^n \frac{\partial f}{\partial x_i}(\mathbf{x}_1, \dots, \mathbf{x}_n) \cdot [-\Delta_i, \Delta_i],$$

where $\tilde{x}_i = (\underline{x}_i + \bar{x}_i)/2$ is the interval's midpoint and $\Delta_i = (\underline{x}_i - \bar{x}_i)/2$ is its half-width.

When all the intervals are the same, e.g., when $\mathbf{x}_i = [0, 1]$, the centered form does not lead to the desired range. Indeed, the centered form always produced an interval centered in the point $f(\tilde{x}_1, \dots, \tilde{x}_n)$. In this case, all midpoints \tilde{x}_i are the same (e.g., equal to 0.5), hence the sample variance $f(\tilde{x}_1, \dots, \tilde{x}_n)$ is equal to 0 on these midpoints. Thus, as a result of applying the centered form, we get an interval centered at 0, i.e., the interval whose lower endpoint is negative. In reality, V is always non-negative, so negative values of V are impossible.

The upper endpoint is also an overestimation: e.g., for $\mathbf{x}_1 = \mathbf{x}_2 = [0, 1]$, we have $\frac{\partial f}{\partial x_1}(x_1, x_2) = x_1 - x_2$, hence

$$\frac{\partial f}{\partial x_1}(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1 - \mathbf{x}_2 = [-1, 1].$$

A similar formula holds for the derivative with respect to x_2 . Since $\Delta_i = 0, 5$, the centered form leads to:

$$f(\mathbf{x}_1, \dots, \mathbf{x}_n) \subseteq 0 + [-1, 1] \cdot [-0.5, 0.5] + [-1, 1] \cdot [-0.5, 0.5] = [-1, 1]$$

– an overestimation in comparison with the actual range $[0, 0.5]$.

1.3 The Problem Reformulated in Statistical Terms

The traditional sample variance is an unbiased estimator for the following problem: observation points x_i satisfy the equation $x_i = u - \varepsilon_i$, where u is an unknown fixed constant and the ε_i are independently and identically distributed random variables with zero expectation and unknown variance σ^2 .

In our paper, we want to handle a situation in which each observation point \tilde{x}_i satisfies the condition $\tilde{x}_i - u - \varepsilon_i \in \Delta_i \cdot [-1, 1]$, where the values Δ_i are assumed to be known. From this model, we can conclude that each $u + \varepsilon_i$ is contained in the corresponding interval $\tilde{x}_i + \Delta_i \cdot [-1, 1] = \mathbf{x}_i$. As a solution to this problem, we take the interval consisting of all the results of applying the estimator V to different values $x_1 \in \mathbf{x}_1, \dots, x_n \in \mathbf{x}_n$.

1.4 Related Statistical Methods Sometimes Overestimate

The above statistical reformulation suggests the possibility of using statistical methods to compute the bounds for variance. If we have n exact observation points, then we can assign, to each of these points, equal probability $1/n$ and thus, form the sample distribution. The cumulative density function (CDF) $F(x)$ for this distribution is easy to describe: we order the observation points in increasing order into a sequence $x_{(1)} \leq \dots \leq x_{(n)}$; then $F(x) = 0$ for $x < x_{(1)}$, $F(x) = 1/n$ for $x \in (x_{(1)}, x_{(2)})$, $F(x) = 2/n$ for $x \in (x_{(2)}, x_{(3)})$, etc., and $F(x) = 1$ for $x > x_{(n)}$. The sample mean and the sample variance are exactly the mean and variance of this sample distribution.

When, instead of the exact observation points, we have the intervals $\mathbf{x}_i = [\underline{x}_i, \bar{x}_i]$, then the CDF $F(x)$ of the sample distribution depends on which values $x_i \in \mathbf{x}_i$ we pick. The smallest possible value $\underline{F}(x)$ of $F(x)$ corresponds to picking $x_i = \bar{x}_i$, and the largest possible value $\overline{F}(x)$ of $F(x)$ corresponds to selecting $x_i = \underline{x}_i$. Thus, for every x , instead of a single value CDF, we have an interval $[\underline{F}(x), \overline{F}(x)]$ of possible values of $F(x)$. This interval-valued CDF is called a *p-bound*, or a *p-box* (see, e.g., [2]).

For each p-box, we can compute the interval of possible values of variance. If we apply this computation to the sample p-box then the resulting interval contains the desired interval for V .

Alas, it is sometimes an overestimation. Indeed, e.g., for a single value $\mathbf{x}_1 = [0, 1]$, the variance is clearly 0, but the p-box is non-trivial:

- $\underline{F}(x) = 0$ for $x < 1$ and $= 1$ for $x \geq 1$;
- $\overline{F}(x) = 0$ for $x < 0$ and $= 1$ for $x \geq 0$.

Within this p-box, there are distributions $F(x) \in [\underline{F}(x), \overline{F}(x)]$ for which the variance is non-zero, hence the resulting variance interval will include some non-zero values – and thus, it will be an overestimation for the desired interval $[0, 0]$.

1.5 We Need New Methods

In short, the existing methods do not always provide us with sharp estimates for \underline{V} and \overline{V} , so we need new methods.

2 First Result: Computing \underline{V}

First, we design a *feasible* algorithm for computing the exact lower bound \underline{V} of the sample variance. Specifically, our algorithm is *quadratic-time*, i.e., it requires $O(n^2)$ computational steps (arithmetic operations or comparisons) for n interval data points $\mathbf{x}_i = [\underline{x}_i, \bar{x}_i]$.

The algorithm \underline{A} is as follows:

- First, we sort all $2n$ values $\underline{x}_i, \bar{x}_i$ into a sequence $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(2n)}$.
- Second, we compute \underline{E} and \bar{E} and select all “small intervals” $[x_{(k)}, x_{(k+1)}]$ that intersect with $[\underline{E}, \bar{E}]$.
- For each of the selected small intervals $[x_{(k)}, x_{(k+1)}]$, we compute the ratio $r_k = S_k/N_k$, where

$$S_k \stackrel{\text{def}}{=} \sum_{i: \underline{x}_i \geq x_{(k+1)}} \underline{x}_i + \sum_{j: \bar{x}_j \leq x_{(k)}} \bar{x}_j,$$

and N_k is the total number of such i 's and j 's. If $r_k \in [x_{(k)}, x_{(k+1)}]$, then we compute

$$V'_k \stackrel{\text{def}}{=} \frac{1}{n-1} \cdot \left(\sum_{i: \underline{x}_i \geq x_{(k+1)}} (\underline{x}_i - r_k)^2 + \sum_{j: \bar{x}_j \leq x_{(k)}} (\bar{x}_j - r_k)^2 \right).$$

If $N_k = 0$, we take $V'_k \stackrel{\text{def}}{=} 0$.

- Finally, we return the smallest of the values V'_k as \underline{V} .

Theorem 2.1. *The algorithm \underline{A} always compute \underline{V} is quadratic time.*

(For readers' convenience, all the proofs are placed in the special Proofs section).

We have implemented this algorithm in C++, it works really fast.

Example. We start with 5 intervals: $\mathbf{x}_1 = [2.1, 2.6]$, $\mathbf{x}_2 = [2.0, 2.1]$, $\mathbf{x}_3 = [2.2, 2.9]$, $\mathbf{x}_4 = [2.5, 2.7]$, and $\mathbf{x}_5 = [2.4, 2.8]$. After sorting the bounds, we get the following “small intervals”: $[x_{(1)}, x_{(2)}] = [2.0, 2.1]$, $[x_{(2)}, x_{(3)}] = [2.1, 2.1]$, $[x_{(3)}, x_{(4)}] = [2.1, 2.2]$, $[x_{(4)}, x_{(5)}] = [2.2, 2.4]$, $[x_{(5)}, x_{(6)}] = [2.4, 2.5]$, $[x_{(6)}, x_{(7)}] = [2.5, 2.6]$, $[x_{(7)}, x_{(8)}] = [2.6, 2.7]$, $[x_{(8)}, x_{(9)}] = [2.7, 2.8]$, and $[x_{(9)}, x_{(10)}] = [2.8, 2.9]$.

The interval for sample average is $\mathbf{E} = [2.24, 2.62]$, so we only keep the following four small intervals that have non-empty intersection with \mathbf{E} : $[x_{(4)}, x_{(5)}] = [2.2, 2.4]$, $[x_{(5)}, x_{(6)}] = [2.4, 2.5]$, $[x_{(6)}, x_{(7)}] = [2.5, 2.6]$, and $[x_{(7)}, x_{(8)}] = [2.6, 2.7]$. For these intervals:

- $S_4 = 7.0$, $N_4 = 3$, so $r_4 = 2.333\dots$;
- $S_5 = 4.6$, $N_5 = 2$, so $r_5 = 2.3$;
- $S_6 = 2.1$, $N_6 = 1$, so $r_6 = 2.1$;
- $S_7 = 4.7$, $N_7 = 2$, so $r_7 = 2.35$.

Of the four values r_k , only r_4 lies within the corresponding small interval. For this small interval, $V'_4 = 0.021666\dots$, so $\underline{V} = 0.021666\dots$

3 Second Result: Computing \overline{V} is NP-Hard

Our second result is that the general problem of computing \overline{V} from given intervals \mathbf{x}_i is computationally difficult, or, in precise terms, NP-hard (for exact definitions of NP-hardness, see, e.g., [4, 5, 8]).

Theorem 3.1. *Computing \overline{V} is NP-hard.*

The very fact that computing the range of a quadratic function is NP-hard was first proven by Vavasis [10] (see also [5]). We have shown that this difficulty happens even for very simple quadratic functions (1.1) frequently used in data processing.

A natural question is: maybe the difficulty comes from the requirement that the range be computed exactly? In practice, it is often sufficient to compute, in a reasonable amount of time, a usefully accurate estimate \tilde{V} for \overline{V} , i.e., an estimate \tilde{V} which is accurate with a given accuracy $\varepsilon > 0$: $|\tilde{V} - \overline{V}| \leq \varepsilon$. Alas, for any ε , such computations are also NP-hard:

Theorem 3.2. *For every $\varepsilon > 0$, the problem of computing \overline{V} with accuracy ε is NP-hard.*

It is worth mentioning that \overline{V} can be computed exactly in exponential time $O(2^n)$:

Theorem 3.3. *There exists an algorithm that computes \overline{V} in exponential time.*

4 Third Result: A Feasible Algorithm That Computes \overline{V} in Many Practical situations

NP-hard means, crudely speaking, that there are no general ways for solving all particular cases of this problem (i.e., computing \overline{V}) in reasonable time.

However, we show that there are algorithms for computing \overline{V} for many reasonable situations. Namely, we propose an efficient algorithm \mathcal{A} that computes \overline{V} for the case when all the interval midpoints (“measured values”) $\tilde{x}_i = (\underline{x}_i + \overline{x}_i)/2$ are definitely different from each other, in the sense that the “narrowed” intervals $[\tilde{x}_i - \Delta_i/n, \tilde{x}_i + \Delta_i/n]$ – where $\Delta_i = (\underline{x}_i - \overline{x}_i)/2$ is the interval’s half-width – do not intersect with each other.

This algorithm $\overline{\mathcal{A}}$ is as follows:

- First, we sort all $2n$ endpoints of the narrowed intervals $\tilde{x}_i - \Delta_i/n$ and $\tilde{x}_i + \Delta_i/n$ into a sequence $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(2n)}$. This enables us to divide the real line into $2n + 2$ segments (“small intervals”) $[x_{(k)}, x_{(k+1)}]$, where we denoted $x_{(0)} \stackrel{\text{def}}{=} -\infty$ and $x_{(2n+1)} \stackrel{\text{def}}{=} +\infty$.

- Second, we compute \underline{E} and \overline{E} and pick all “small intervals” $[x_{(k)}, x_{(k+1)}]$ that intersect with $[\underline{E}, \overline{E}]$.
- For each of remaining small intervals $[x_{(k)}, x_{(k+1)}]$, for each i from 1 to n , we pick the following value of x_i :
 - if $x_{(k+1)} < \tilde{x}_i - \Delta_i/n$, then we pick $x_i = \overline{x}_i$;
 - if $x_{(k)} > \tilde{x}_i + \Delta_i/n$, then we pick $x_i = \underline{x}_i$;
 - for all other i , we consider both possible values $x_i = \overline{x}_i$ and $x_i = \underline{x}_i$.

As a result, we get one or several sequences of x_i . For each of these sequences, we check whether the average E of the selected values x_1, \dots, x_n is indeed within this small interval, and if it is, compute the sample variance by using the formula (1.1).

- Finally, we return the largest of the computed sample variances as \overline{V} .

Theorem 4.1. *The algorithm \overline{A} compute \overline{V} is quadratic time for all the cases in which the “narrowed” intervals do not intersect with each other.*

This algorithm also works when, for some fixed k , no more than k “narrowed” intervals can have a common point:

Theorem 4.2. *For every positive integer k , the algorithm \overline{A} compute \overline{V} is quadratic time for all the cases in which no more than k “narrowed” intervals can have a common point.*

This computation time is quadratic in n but it grows exponentially with k . So, when k grows, this algorithm requires more and more computation time. It is worth mentioning that the examples on which we prove NP-hardness (see proof of Theorem 3.1) correspond to the case when all n narrowed intervals have a common point.

5 Sample Mean, Sample Variance: What Next?

5.1 Sample Covariance

When we have two sets of data x_1, \dots, x_n and y_1, \dots, y_n , we normally compute *sample covariance*

$$C = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y}),$$

where

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i; \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i.$$

Sample covariance is used to describe the correlation between x_i and y_i . If we take interval uncertainty into consideration, then, after each measurement, we do not get the exact values of $x_1, \dots, x_n, y_1, \dots, y_n$; instead, we only have *intervals* $[\underline{x}_1, \bar{x}_1], \dots, [\underline{x}_n, \bar{x}_n], [\underline{y}_1, \bar{y}_1], \dots, [\underline{y}_n, \bar{y}_n]$. Depending on what are the actual values of $x_1, \dots, x_n, y_1, \dots, y_n$ within these intervals, we get different values of sample covariance. To take the interval uncertainty into consideration, we need to be able to describe the interval $[\underline{C}, \bar{C}]$ of possible values of the sample covariance C .

So, we arrive at the following problems: given the intervals $[\underline{x}_i, \bar{x}_i], [\underline{y}_i, \bar{y}_i]$, compute the lower and upper bounds \underline{C} and \bar{C} for the interval of possible values of sample covariance.

It turns out that these problems are also NP-hard:

Theorem 5.1. *The problem of computing \bar{C} from the interval inputs $[\underline{x}_i, \bar{x}_i], [\underline{y}_i, \bar{y}_i]$ is NP-hard.*

Theorem 5.2. *The problem of computing \underline{C} from the interval inputs $[\underline{x}_i, \bar{x}_i], [\underline{y}_i, \bar{y}_i]$ is NP-hard.*

Comment. These results were first announced in [7].

5.2 Sample Correlation

As we have mentioned, sample covariance C between the data sets x_1, \dots, x_n and y_1, \dots, y_n is often used to compute sample correlation

$$\rho = \frac{C}{\sigma_x \cdot \sigma_y}, \quad (5.1)$$

where $\sigma_x = \sqrt{V_x}$ is the sample standard deviation of the values x_1, \dots, x_n , and $\sigma_y = \sqrt{V_y}$ is the sample standard deviation of the values y_1, \dots, y_n .

When we only have *intervals* $[\underline{x}_1, \bar{x}_1], \dots, [\underline{x}_n, \bar{x}_n], [\underline{y}_1, \bar{y}_1], \dots, [\underline{y}_n, \bar{y}_n]$, we have an interval $[\underline{\rho}, \bar{\rho}]$ of possible value of correlation. It turns out that, similar to sample covariance, computation of the endpoints of this interval problems is also an NP-hard problem:

Theorem 5.3. *The problem of computing $\bar{\rho}$ from the interval inputs $[\underline{x}_i, \bar{x}_i], [\underline{y}_i, \bar{y}_i]$ is NP-hard.*

Theorem 5.4. *The problem of computing $\underline{\rho}$ from the interval inputs $[\underline{x}_i, \bar{x}_i], [\underline{y}_i, \bar{y}_i]$ is NP-hard.*

Comment. The fact that the problems of computing sample covariance and sample correlation are NP-hard means that, crudely speaking, that there is no feasible algorithm that would always compute the desired bounds for C and ρ . A similar NP-hardness result holds for sample variance, but in that case, we were also able to produce a feasible algorithm that works in many practical cases. It is desirable to design similar algorithms for sample covariance and sample correlation.

5.3 Other Characteristics

Not all statistical characteristics are difficult to compute for interval data, some are easy. In addition to sample mean, we can mention sample *median*. Since the median is increasing in x_1, \dots, x_n , its smallest possible value is attained for $\underline{x}_1, \dots, \underline{x}_n$, and its largest possible value is attained for $\bar{x}_1, \dots, \bar{x}_n$.

It is desirable to analyze other statistical characteristics from this viewpoint.

6 Proofs

Proof of Theorem 2.1

1°. Let us first show that the algorithm described in Section 2 is indeed correct.

1.1°. Indeed, let $x_1^{(0)} \in \mathbf{x}_1, \dots, x_n^{(0)} \in \mathbf{x}_n$ be the values for which the sample variance V attains minimum on the box $\mathbf{x}_1 \times \dots \times \mathbf{x}_n$.

Let us pick one of the n variables x_i , and let fix the values of all the other variables x_j ($j \neq i$) at $x_j = x_j^{(0)}$. When we substitute $x_j = x_j^{(0)}$ for all $j \neq i$ into the expression for sample variance, V becomes a quadratic function of x_i .

This function of one variable should attain its minimum on the interval \mathbf{x}_i at the value $x_i^{(0)}$.

1.2°. Let us start with the analysis of the quadratic function of one variable we described in Part 1.1 of this proof.

By definition, the sample variance V is a sum of non-negative terms; thus, its value is always non-negative. Therefore, the corresponding quadratic function of one variable always has a global minimum. This function is decreasing before this global minimum, and increasing after it.

1.3°. Where is the global minimum of the quadratic function of one variable described in Part 1.1?

It is attained when $\partial V / \partial x_i = 0$. Differentiating the formula (1.1) with respect to x_i , we conclude that

$$\frac{\partial V}{\partial x_i} = \frac{1}{n-1} \cdot \left(2(x_i - E) + \sum_{j=1}^n 2(E - x_j) \cdot \frac{\partial E}{\partial x_j} \right). \quad (6.1)$$

Since $\partial E / \partial x_i = 1/n$, we conclude that

$$\frac{\partial V}{\partial x_i} = \frac{2}{n-1} \cdot \left((x_i - E) + \sum_{j=1}^n (E - x_j) \cdot \frac{1}{n} \right). \quad (6.2)$$

Here,

$$\sum_{j=1}^n (E - x_j) = n \cdot E - \sum_{j=1}^n x_j. \quad (6.3)$$

By definition of the average E , this difference is 0, hence the formula (6.2) takes the form

$$\frac{\partial V}{\partial x_i} = \frac{2}{n-1} \cdot (x_i - E).$$

So, this function attains the minimum when $x_i - E = 0$, i.e., when $x_i = E$.

Since

$$E = \frac{x_i}{n} + \frac{\sum'_i x_j}{n},$$

where \sum'_i means the sum over all $j \neq i$, the equality $x_i = E$ means that

$$x_i = \frac{x_i}{n} + \frac{\sum'_i x_j^{(0)}}{n}.$$

Moving terms containing x_i into the left-hand side and dividing by the coefficient at x_i , we conclude that the minimum is attained when

$$x_i = E'_i \stackrel{\text{def}}{=} \frac{\sum'_i x_j^{(0)}}{n-1},$$

i.e., when x_i is equal to the arithmetic average E'_i of all other elements.

1.4°. Let us now use the knowledge of a global minimum to describe where the desired function attains its minimum on the interval \mathbf{x}_i .

In our general description of non-negative quadratic functions of one variable, we mentioned that each such function is decreasing before the global minimum and increasing after it. Thus, for $x_i < E'_i$, the function V is decreasing, for $x_i > E'_i$, this function is increasing. Therefore:

- If $E'_i \in \mathbf{x}_i$, the global minimum of the function V of one variable is attained within the interval \mathbf{x}_i , hence the minimum on the interval \mathbf{x}_i is attained for $x_i = E'_i$.
- If $E'_i < \underline{x}_i$, the function V is increasing on the interval \mathbf{x}_i and therefore, its minimum on this interval is attained when $x_i = \underline{x}_i$.
- Finally, if $E'_i > \bar{x}_i$, the function V is decreasing on the interval \mathbf{x}_i and therefore, its minimum on this interval is attained when $x_i = \bar{x}_i$.

1.5°. Let us reformulate the above conditions in terms of the average

$$E = \frac{1}{n} \cdot x_i + \frac{n-1}{n} \cdot E'_i.$$

- In the first case, when $x_i = E'_i$, we have $x_i = E = E'_i$, so $E \in \mathbf{x}_i$.
- In the second case, we have $E'_i < \underline{x}_i$ and $x_i = \underline{x}_i$. Therefore, in this case, $E < \underline{x}_i$.
- In the third case, we have $E'_i > \bar{x}_i$ and $x_i = \bar{x}_i$. Therefore, in this case, $E > \bar{x}_i$.

Thus:

- If $E \in \mathbf{x}_i$, then we cannot be in the second or third cases. Thus, we are in the first case, hence $x_i = E$.
- If $E < \underline{x}_i$, then we cannot be in the first or the third cases. Thus, we are the second case, hence $x_i = \underline{x}_i$.
- If $E > \bar{x}_i$, then we cannot be in the first or the second cases. Thus, we are in the third case, hence $x_i = \bar{x}_i$.

1.6°. So, as soon as we determine the position of E with respect to all the bounds \underline{x}_i and \bar{x}_i , we will have a pretty good understanding of all the values x_i at which the minimum is attained.

Hence, to find the minimum, we will analyze how the endpoints \underline{x}_i and \bar{x}_i divide the real line, and consider all the resulting sub-intervals.

Let the corresponding subinterval $[x_{(k)}, x_{(k+1)}]$ be fixed. For the i 's for which $E \notin \mathbf{x}_i$, the values x_i that correspond to the minimal sample variance are uniquely determined by the above formulas.

For the i 's for which $E \in \mathbf{x}_i$ the selected value x_i should be equal to E . To determine this E , we can use the fact that E is equal to the average of all thus selected values x_i , in other words, that we should have

$$E = \frac{1}{n} \cdot \left(\sum_{i: \underline{x}_i \geq x_{(k+1)}} \underline{x}_i + (n - N_k) \cdot E + \sum_{j: \bar{x}_j \leq x_{(k)}} \bar{x}_j \right), \quad (6.4)$$

where $(n - N_k) \cdot E$ combines all the points for which $E \in \mathbf{x}_i$. Multiplying both sides of (6.4) by n and subtracting $n \cdot E$ from both sides, we conclude that (in notations of Section 2), we have $E = S_k/N_k$ – what we denoted, in the algorithm's description, by r_k . If thus defined r_k does not belong to the subinterval $[x_{(k)}, x_{(k+1)}]$, this contradiction with our initial assumption shows that there cannot be any minimum in this subinterval, so this subinterval can be easily dismissed.

The corresponding sample variance is denoted by V'_k . If $N_k = 0$, this means that E belongs to all the intervals \mathbf{x}_i and therefore, that the lower endpoint \underline{V} is exactly 0 – so we assign $V'_k = 0$.

2°. To complete the proof of Theorem 2.1, we must show that this algorithm indeed requires quadratic time.

Indeed, sorting requires $O(n \cdot \log(n))$ steps (see, e.g., [1]), and the rest of the algorithm requires linear time ($O(n)$) for each of $2n$ subintervals, i.e., the total quadratic time.

The theorem is proven.

Proof of Theorem 3.1

1°. By definition, a problem is NP-hard if any problem from the class NP can be reduced to it. Therefore, to prove that a problem \mathcal{P} is NP-hard, it is sufficient to reduce one of the known NP-hard problems \mathcal{P}_0 to \mathcal{P} .

In this case, since \mathcal{P}_0 is known to be NP-hard, this means that every problem from the class NP can be reduced to \mathcal{P}_0 , and since \mathcal{P}_0 can be reduced to \mathcal{P} , thus, the original problem from the class NP is reducible to \mathcal{P} .

For our proof, as the known NP-hard problem \mathcal{P}_0 , we take a *subset* problem: given n positive integers s_1, \dots, s_n , to check whether there exist signs $\eta_i \in \{-1, +1\}$ for which the signed sum $\sum_{i=1}^n \eta_i \cdot s_i$ equals 0.

We will show that this problem can be reduced to the problem of computing \overline{V} , i.e., that to every instance (s_1, \dots, s_n) of the problem \mathcal{P}_0 , we can put into correspondence such an instance of the \overline{C} -computing problem that based on its solution, we can easily check whether the desired signs exist.

As this instance, we take the instance corresponding to the intervals $[\underline{x}_i, \overline{x}_i] = [-s_i, s_i]$. We want to show that for the corresponding problem, $\overline{V} = C_0$, where we denoted

$$C_0 \stackrel{\text{def}}{=} \frac{1}{n-1} \cdot \sum_{i=1}^n s_i^2, \quad (6.5)$$

if and only if there exist signs η_i for which $\sum \eta_i \cdot s_i = 0$.

2°. Let us first show that in all cases, $\overline{V} \leq C_0$.

Indeed, it is known that the formula for the sample variance can be reformulated in the following equivalent form:

$$V = \frac{1}{n-1} \cdot \sum_{i=1}^n x_i^2 - \frac{n}{n-1} \cdot (\overline{x})^2. \quad (6.6)$$

Since $x_i \in [-s_i, s_i]$, we can conclude that $x_i^2 \leq s_i^2$ hence $\sum x_i^2 \leq \sum s_i^2$. Since $(\overline{x})^2 \geq 0$, we thus conclude that

$$V \leq \frac{1}{n-1} \cdot \sum_{i=1}^n s_i^2 = C_0.$$

In other words, every possible value V of the sample variance is smaller than or equal to C_0 . Thus, the largest of these possible values, i.e., \overline{V} , also cannot exceed C_0 , i.e., $\overline{V} \leq C_0$.

3°. Let us now prove that if the desired signs η_i exist, then $\bar{V} = C_0$.

Indeed, in this case, for $x_i = \eta_i \cdot s_i$, we have $\bar{x} = 0$ and $x_i^2 = s_i^2$, hence

$$V = \frac{1}{n-1} \cdot \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \cdot \sum_{i=1}^n s_i^2 = C_0.$$

So, the sample variance V is always $\leq C_0$, and it attains the value C_0 for some x_i . Therefore, $\bar{V} = C_0$.

4°. To complete the proof of Theorem 3.1, we must show that, vice versa, if $\bar{V} = C_0$, then the desired signs exist.

Indeed, let $\bar{V} = C_0$. Sample variance is a continuous function on a compact set $\mathbf{x}_1 \times \dots \times \mathbf{x}_n$, hence its maximum on this compact set is attained for some values $x_1 \in \mathbf{x}_1 = [-s_1, s_1], \dots, x_n \in \mathbf{x}_n = [-s_n, s_n]$. In other words, for the corresponding values of x_i , the sample variance V is equal to C_0 .

Since $x_i \in [-s_i, s_i]$, we can conclude that $x_i^2 \leq s_i^2$; since $(\bar{x})^2 \geq 0$, we get $V \leq C_0$. If $|x_i|^2 < s_i^2$ or $(\bar{x})^2 > 0$, then we would have $V < C_0$. Thus, the only way to have $V = C_0$ is to have $x_i^2 = s_i^2$ and $\bar{x} = 0$. The first equality leads to $x_i = \pm s_i$, i.e., to $x_i = \eta_i \cdot s_i$ for some $\eta_i \in \{-1, +1\}$. Since \bar{x} is, by definition, the (arithmetic) average of the values x_i , the equality $\bar{x} = 0$ then

leads to $\sum_{i=1}^n \eta_i \cdot s_i = 0$. So, if $\bar{V} = C_0$, then the desired signs do exist.

The theorem is proven.

Proof of Theorem 3.2

1°. Let $\varepsilon > 0$ be fixed. We will show that the subset problem can be reduced to the problem of computing \bar{V} with accuracy ε , i.e., that to every instance (s_1, \dots, s_n) of the subset problem \mathcal{P}_0 , we can put into correspondence such an instance of the ε -approximate \bar{C} -computation problem that based on its solution, we can easily check whether the desired signs exist.

For this reduction, we will use two parameters. The first one – C_0 – is the same as in the proof of Theorem 3.1. We will also need a new real-valued parameter k ; its value depend on ε and n . We could produce this value right away, but we believe that the proof will be much clearer if we keep it undetermined until it becomes clear what value k we need to choose for the proof to be valid.

As the desired instance, we take the instance corresponding to the intervals $[\underline{x}_i, \bar{x}_i] = [-k \cdot s_i, k \cdot s_i]$ for an appropriate value k . Let \tilde{V} be a number produced, for this problem, by a ε -accurate computation algorithm, i.e., a number for which $|\tilde{V} - \bar{V}| \leq \varepsilon$. We want to show that $\tilde{V} \geq k^2 \cdot C_0 - \varepsilon$ if and only if there exist signs η_i for which $\sum \eta_i \cdot s_i = 0$.

2°. When we multiply each value x_i by a constant k , the sample variance is multiplied by k^2 . As a result, the upper bound \bar{V} corresponding to $x_i \in$

$[-k \cdot s_i, k \cdot s_i]$ is exactly k^2 times larger than the upper bound \bar{v} corresponding to k times smaller values $z_i \in [-s_i, s_i]$: $\bar{v} = \bar{V}/k^2$.

Hence, when \tilde{V} approximates \bar{V} with an accuracy ε , the corresponding value $\tilde{v} \stackrel{\text{def}}{=} \tilde{V}/k^2$ approximates $\bar{v} (= \bar{V}/k^2)$ with the accuracy $\delta \stackrel{\text{def}}{=} \varepsilon/k^2$.

In terms of \tilde{v} , the above inequality $\tilde{V} \geq k^2 \cdot C_0 - \varepsilon$ takes the following equivalent form: $\tilde{v} \geq C_0 - \delta$.

Thus, in terms of \tilde{v} , the desired property can be formulated as follows: $\tilde{v} \geq C_0 - \delta$ if and only if there exist signs η_i for which $\sum \eta_i \cdot s_i = 0$.

3°. Let us first show that if the desired signs η_i exist, then $\tilde{v} \geq C_0 - \delta$.

Indeed, in this case, similarly to the proof of Theorem 3.1, we can conclude that $\bar{v} = C_0$. Since \tilde{v} is a δ -approximation to the actual upper bound \bar{v} , we can therefore conclude that $\tilde{v} \geq \bar{v} - \delta = C_0 - \delta$. The statement is proven.

4°. Vice versa, let us assume that $\tilde{v} \geq C_0 - \delta$. Let us prove that in this case, the desired signs exist.

4.1°. Since \tilde{v} is a δ -approximation to the upper bound \bar{v} , we thus conclude that $\bar{v} \geq \tilde{v} - \delta$ and therefore, $\bar{v} \geq C_0 - 2\delta$.

Similarly to the proof of Theorem 3.1, we can conclude that the maximum is attained for some values $z_i \in [-s_i, s_i]$ and therefore, there exist values $z_i \in [-s_i, s_i]$ for which the sample variance v exceeds $C_0 - 2\delta$:

$$v \stackrel{\text{def}}{=} \frac{1}{n-1} \cdot \sum_{i=1}^n z_i^2 - \frac{n}{n-1} \cdot (\bar{z})^2 \geq C_0 - 2\delta,$$

i.e., substituting the expression (6.5) for C_0 , that

$$\frac{1}{n-1} \cdot \sum_{i=1}^n z_i^2 - \frac{n}{n-1} \cdot (\bar{z})^2 \geq \frac{1}{n-1} \cdot \sum_{i=1}^n s_i^2 - 2\delta. \quad (6.7)$$

4.2°. The following proof will be similar to the corresponding part of the proof of Theorem 3.1. The main difference is that we have approximate equalities instead of exact ones:

- In the proof of Theorem 3.1, we used the fact that $V = C_0$ to prove that the corresponding values x_i are equal to $\pm s_i$, and that their sum is equal to 0.
- Here, v is only approximately equal to C_0 . As a result, we will only be able to show that the values z_i are *close* to $\pm s_i$, and that the sum of z_i is *close* to 0. From these closenesses, we will then be able to conclude (for sufficiently large k) that the sum of the corresponding terms $\pm s_i$ is exactly equal to 0.

4.3°. Let us first prove that for every i , the value z_i^2 is close to s_i^2 . Specifically, we know that $z_i^2 \leq s_i^2$; we will prove that

$$z_i^2 \geq s_i^2 - 2(n-1) \cdot \delta. \quad (6.8)$$

We will prove this inequality by reduction to a contradiction. Indeed, let us assume that for some i_0 , this inequality is not true. This means that

$$z_{i_0}^2 < s_{i_0}^2 - 2(n-1) \cdot \delta. \quad (6.9)$$

Since $z_i \in [-s_i, s_i]$, for all i , in particular, for all $i \neq i_0$, we conclude, for all $i \neq i_0$, that

$$z_i^2 \leq s_i^2. \quad (6.10)$$

Adding the inequality (6.9) and $(n-1)$ inequalities (6.10) corresponding to all values $i \neq i_0$, we get

$$\sum_{i=1}^n z_i^2 < \sum_{i=1}^n s_i^2 - 2(n-1) \cdot \delta. \quad (6.11)$$

Dividing both sides of this inequality by $n-1$, we get a contradiction with (6.7). This contradiction shows that (6.8) indeed holds for every i .

4.4°. The inequality (6.8) says, crudely speaking, that z_i^2 is close to s_i^2 . According to our “action plan” (as outlined in Part 4.2 of this proof), we want to conclude that z_i is close to $\pm s_i$, i.e., that $|z_i|$ is close to s_i .

To be able to make a meaningful conclusion about z_i from the inequality (6.8), we must make sure that the right-hand side of the inequality (6.8) is positive: otherwise, this inequality is true simply because its left-hand side is non-negative, and the right-hand side is non-positive.

The value s_i is a positive integer, so $s_i^2 \geq 1$. Therefore, to guarantee that the right-hand side of (6.8) is positive, it is sufficient to select k for which, for the corresponding value $\delta = \varepsilon/k^2$, we have

$$2(n-1) \cdot \delta < 1. \quad (6.12)$$

In the following text, we will assume that this condition is indeed satisfied.

4.5°. Let us show that under the condition (6.12), the value $|z_i|$ is indeed close to s_i . To be more precise, we already know that $|z_i| \leq s_i$; we are going to prove that

$$|z_i| \geq s_i - 2(n-1) \cdot \delta. \quad (6.13)$$

Indeed, since the right-hand side of the inequality (6.8) is supposed to be close to s_i , it makes sense to represent it as s_i^2 times a factor close to 1. To be more precise, we reformulate the inequality (6.8) in the following equivalent form:

$$z_i^2 \geq s_i^2 \cdot \left(1 - \frac{2(n-1) \cdot \delta}{s_i^2}\right). \quad (6.14)$$

Since both sides of this inequality are non-negative, we can extract the square root from both sides and get the following inequality:

$$|z_i| \geq s_i \cdot \sqrt{1 - \frac{2(n-1) \cdot \delta}{s_i^2}}. \quad (6.15)$$

The square root in the right-hand side of (6.15) is of the type $\sqrt{1-t}$, with $0 \leq t \leq 1$. It is known that for such t , we have $\sqrt{1-t} \geq 1-t$. Therefore, from (6.15), we can conclude that

$$|z_i| \geq s_i \cdot \sqrt{1 - \frac{2(n-1) \cdot \delta}{s_i^2}} \geq s_i \cdot \left(1 - \frac{2(n-1) \cdot \delta}{s_i^2}\right),$$

i.e., that

$$|z_i| \geq s_i - \frac{2(n-1) \cdot \delta}{s_i}.$$

Since $s_i \geq 1$, we have

$$\frac{2(n-1) \cdot \delta}{s_i} \leq 2(n-1) \cdot \delta,$$

hence

$$|z_i| \geq s_i - \frac{2(n-1) \cdot \delta}{s_i} \geq s_i - 2(n-1) \cdot \delta.$$

So, the inequality (6.13) is proven.

4.6°. Let us now prove that for the values z_i selected on Step 4.1, the average \bar{z} is close to 0. To be more precise, we will prove that

$$(\bar{z})^2 \leq \frac{n-1}{n} \cdot 2\delta. \quad (6.16)$$

Similarly to Part 4.3 of this proof, we will prove this inequality by reduction to a contradiction. Indeed, assume that this inequality is not true, i.e., that

$$(\bar{z})^2 > \frac{n-1}{n} \cdot 2\delta. \quad (6.17)$$

Since $z_i^2 \leq s_i^2$, we therefore conclude that

$$\sum_{i=1}^n z_i^2 \leq \sum_{i=1}^n s_i^2,$$

hence

$$\frac{1}{n-1} \cdot \sum_{i=1}^n z_i^2 \leq \frac{1}{n-1} \cdot \sum_{i=1}^n s_i^2. \quad (6.18)$$

Adding, to both sides of the inequality (6.18), the inequality (6.17) multiplied by $n/(n-1)$, we get an inequality

$$\frac{1}{n-1} \cdot \sum_{i=1}^n z_i^2 - \frac{n}{n-1} \cdot (\bar{z})^2 < \frac{1}{n-1} \sum_{i=1}^n s_i^2 - 2\delta,$$

which contradicts to (6.7). This contradiction proves that that the inequality (6.16) is true.

4.7°. From the fact that the average \bar{z} is close to 0, we can now conclude that the sum $\sum z_i$ is also close to 0. Specifically, we will now prove that

$$\left| \sum_{i=1}^n z_i \right| \leq n \cdot \sqrt{2\delta}. \quad (6.19)$$

Indeed, from (6.16), we conclude that $(\bar{z})^2 \leq 2\delta$, hence $|\bar{z}| \leq \sqrt{2\delta}$. Multiplying both sides of this inequality by n , we get the desired inequality (6.19).

4.8°. Let us now show that for appropriately chosen k , we will be able to conclude that there exist signs η_i for which $\sum \eta_i \cdot s_i = 0$.

From the inequalities (6.13) and $|z_i| \leq s_i$, we conclude that

$$|s_i - |z_i|| \leq 2(n-1) \cdot \delta. \quad (6.20)$$

Hence, $|z_i| \leq s_i - 2(n-1) \cdot \delta$. Each value s_i is a positive integer, so $s_i \geq 1$. Due to the inequality (6.12), we have $2(n-1) \cdot \delta < 1$, so $|z_i| > 1 - 1 = 0$. Therefore, $z_i \neq 0$, hence each value z_i has a sign. Let us take, as η_i , the sign of the value z_i . Then, the inequality (6.20) takes the form

$$|\eta_i \cdot s_i - z_i| \leq 2(n-1) \cdot \delta. \quad (6.21)$$

Since the absolute value of the sum cannot exceed the sum of absolute values, we therefore conclude that

$$\begin{aligned} \left| \sum_{i=1}^n \eta_i \cdot s_i - \sum_{i=1}^n z_i \right| &= \left| \sum_{i=1}^n (\eta_i \cdot s_i - z_i) \right| \leq \sum_{i=1}^n |\eta_i \cdot s_i - z_i| \leq \\ &\sum_{i=1}^n 2(n-1) \cdot \delta = 2n \cdot (n-1) \cdot \delta. \end{aligned} \quad (6.22)$$

From (6.22) and (6.19), we conclude that

$$\left| \sum_{i=1}^n \eta_i \cdot s_i \right| \leq \left| \sum_{i=1}^n z_i \right| + \left| \sum_{i=1}^n \eta_i \cdot s_i - \sum_{i=1}^n z_i \right| = n \cdot \sqrt{2\delta} + 2n \cdot (n-1) \cdot \delta. \quad (6.23)$$

All values s_i are integers, hence, the sum $\sum \eta_i \cdot s_i$ is also an integer, and so is its absolute value $|\sum \eta_i \cdot s_i|$. Thus, if we select k for which the right-hand side of the inequality (6.23) is less than 1, i.e., for which

$$n \cdot \sqrt{2\delta} + 2n \cdot (n-1) \cdot \delta < 1, \quad (6.24)$$

we therefore conclude that the absolute value of an integer $\sum \eta_i \cdot s_i$ is smaller than 1, so it must be equal to 0: $\sum \eta_i \cdot s_i = 0$.

Thus, to complete the proof, it is sufficient to find k for which, for the corresponding value $\delta = \varepsilon/k^2$, both the inequalities (6.12) and (6.24) hold. To guarantee the inequality (6.24), it is sufficient to have

$$n \cdot \sqrt{2\delta} \leq \frac{1}{3} \quad (6.25)$$

and

$$2n \cdot (n-1) \cdot \delta \leq \frac{1}{3}. \quad (6.26)$$

The inequality (6.25) is equivalent to

$$\delta \leq \frac{1}{18n^2};$$

the inequality (6.26) is equivalent to

$$\delta \leq \frac{1}{6n \cdot (n-1)};$$

and the inequality (6.12) is equivalent to

$$\delta \leq \frac{1}{2(n-1)}.$$

Thus, to satisfy all three inequalities, we must choose δ for which $\delta = \varepsilon/k^2 = \delta_0$, where we denoted

$$\delta_0 \stackrel{\text{def}}{=} \min \left(\frac{1}{18n^2}, \frac{1}{6n \cdot (n-1)}, \frac{1}{2(n-1)} \right).$$

The original expression (1.1) for the sample variance only works for $n \geq 2$. For such n , $18n^2 > 6n \cdot (n-1)$ and $18n^2 > 2(n-1)$, hence the above formula can be simplified into

$$\delta_0 = \frac{1}{18n^2}.$$

To get this δ as $\delta_0 = \varepsilon/k^2$, we must take $k = \sqrt{\varepsilon/\delta_0} = 3n \cdot \sqrt{2\varepsilon}$. For this k , as we have shown before, the reduction holds, so the theorem is proven.

Proof of Theorem 3.3

Let $x_1^{(0)} \in \mathbf{x}_1, \dots, x_n^{(0)} \in \mathbf{x}_n$ be the values for which the sample variance V attains maximum on the box $\mathbf{x}_1 \times \dots \times \mathbf{x}_n$.

Let us pick one of the n variables x_i , and let fix the values of all the other variables x_j ($j \neq i$) at $x_j = x_j^{(0)}$. When we substitute $x_j = x_j^{(0)}$ for all $j \neq i$ into the expression for sample variance, V becomes a quadratic function of x_i .

This function of one variable should attain its maximum on the interval \mathbf{x}_i at the value $x_i^{(0)}$.

As we have mentioned in the proof of Theorem 2.1, by definition, the sample variance V is a sum of non-negative terms; thus, its value is always non-negative. Therefore, the corresponding quadratic function of one variable always has a global minimum. This function is decreasing before this global minimum, and increasing after it. Thus, its maximum on the interval \mathbf{x}_i is attained at one of the endpoints of this interval.

In other words, for each variable x_i , the maximum is attained either for $x_i = \underline{x}_i$, or for $x_i = \bar{x}_i$. Thus, to find \bar{V} , it is sufficient to compute V for 2^n possible combinations $(x_1^\pm, \dots, x_n^\pm)$, where $x_i^- \stackrel{\text{def}}{=} \underline{x}_i$ and $x_i^+ \stackrel{\text{def}}{=} \bar{x}_i$, and find the largest of the resulting 2^n numbers.

Proof of Theorems 4.1 and 4.2

1°. Similarly to the proof of Theorem 2.1, let us first show that the algorithm described in Section 4 is indeed correct.

2°. Similarly to the proof of Theorem 2.1, let x_1, \dots, x_n be the values at which the sample variance attain its maximum on the box $\mathbf{x}_1 \times \dots \times \mathbf{x}_n$. If we fix the values of all the variables but one x_i , then V becomes a quadratic function of x_i . When the function V attains maximum over $x_1 \in \mathbf{x}_1, \dots, x_n \in \mathbf{x}_n$, then this quadratic function of one variable will attain its maximum on the interval \mathbf{x}_i at the point x_i .

We have already shown, in the proof of Theorem 2.1, that this quadratic function has a (global) minimum at $x_i = E'_i$, where E'_i is the average of all the values x_1, \dots, x_n except for x_i . Since this quadratic function of one variable is always non-negative, it cannot have a global maximum. Therefore, its maximum on the interval $\mathbf{x}_i = [\underline{x}_i, \bar{x}_i]$ is attained at one of the endpoints of this interval.

An arbitrary quadratic function of one variable is symmetric with respect to the location of its global minimum, so its maximum on any interval is attained at the point which is the farthest from the minimum. There is exactly one point which is equally close to both endpoints of the interval \mathbf{x}_i : its midpoint \tilde{x}_i . Depending on whether the global minimum is to the left, to the right, or exactly at the midpoint, we get the following three possible cases:

1. If the global minimum E'_i is to the left of the midpoint \tilde{x}_i , i.e., if $E'_i < \tilde{x}_i$, then the upper endpoint is the farthest from E'_i . In this case, the maximum

of the quadratic function is attained at its upper endpoint, i.e., $x_i = \bar{x}_i$.

2. Similarly, if the global minimum E'_i is to the right of the midpoint \tilde{x}_i , i.e., if $E'_i > \tilde{x}_i$, then the lower endpoint is the farthest from E'_i . In this case, the maximum of the quadratic function is attained at its lower endpoint, i.e., $x_i = \underline{x}_i$.
3. If $E'_i = \tilde{x}_i$, then the maximum of V is attained at both endpoints of the interval $\mathbf{x}_i = [\underline{x}_i, \bar{x}_i]$.

3°. In the third case, we have either $x_i = \underline{x}_i$ or $x_i = \bar{x}_i$. Depending on whether x_i is equal to the lower or to the upper endpoints, we can “combine” the corresponding situations with Cases 1 and 2. As a result, we arrive at the conclusion that one of the following two situations happen:

1. either $E'_i \leq \tilde{x}_i$ and $x_i = \bar{x}_i$;
2. either $E'_i \geq \tilde{x}_i$ and $x_i = \underline{x}_i$.

4°. Similarly to the proof of Theorem 2.1, let us reformulate these conclusions in terms of the average E of the maximizing values x_1, \dots, x_n .

The average E'_i can be described as

$$\frac{\sum'_i x_j}{n-1},$$

where \sum'_i means the sum over all $j \neq i$. By definition, $\sum'_j x_j = \sum_j x_j - x_i$, where $\sum_j x_j$ means the sum over all possible j . By definition of E , we have

$$E = \frac{\sum_j x_j}{n},$$

hence $\sum_j x_j = n \cdot E$. Therefore,

$$E'_i = \frac{n \cdot E - x_i}{n-1}.$$

Let us apply this formula to the above three cases.

4.1°. In the first case, we have $\tilde{x}_i \geq E'_i$. So, in terms of E , we get the inequality

$$\tilde{x}_i \geq \frac{n \cdot E - x_i}{n-1}.$$

Multiplying both sides of this inequality by $n-1$, and using the fact that in this case, $x_i = \bar{x}_i = \tilde{x}_i + \Delta_i$, we conclude that

$$(n-1) \cdot \tilde{x}_i \geq n \cdot E - \tilde{x}_i - \Delta_i.$$

Moving all the terms but $n \cdot E$ to the left-hand side and dividing by E , we get the following inequality:

$$E \leq \tilde{x}_i + \frac{\Delta_i}{n}.$$

4.2°. In the second case, we have $\tilde{x}_i \leq E'_i$. So, in terms of E , we get the inequality

$$\tilde{x}_i \leq \frac{n \cdot E - x_i}{n - 1}.$$

Multiplying both sides of this inequality by $n - 1$, and using the fact that in this case, $x_i = \underline{x}_i = \tilde{x}_i - \Delta_i$, we conclude that

$$(n - 1) \cdot \tilde{x}_i \leq n \cdot E - \tilde{x}_i + \Delta_i.$$

Moving all the terms but $n \cdot E$ to the left-hand side and dividing by E , we get the following inequality:

$$E \geq \tilde{x}_i - \frac{\Delta_i}{n}.$$

5°. Parts 4.1 and 4.2 of this proof can be summarized as follows:

- In Case 1, we have $E \leq \tilde{x}_i + \Delta_i/n$ and $x_i = \bar{x}_i$.
- In Case 2, we have $E \geq \tilde{x}_i - \Delta_i/n$ and $x_i = \underline{x}_i$.

Therefore:

- If $E < \tilde{x}_i - \Delta_i/n$, this means that we cannot be in Case 2. So we must be in Case 1 and therefore, we must have $x_i = \bar{x}_i$.
- If $E > \tilde{x}_i + \Delta_i/n$, this means that we cannot be in Case 1. So, we must be in Case 2 and therefore, we must have $x_i = \underline{x}_i$.

The only case when we do not know which endpoint for x_i we should choose is the case when E belongs to the narrowed interval $[\tilde{x}_i - \Delta/n, \tilde{x}_i + \Delta_i]$.

6°. Hence, once we know where E is with respect to the endpoints of all narrowed intervals, we can determine the values of all optimal x_i – except for those that are within this narrowed interval. Since we consider the case when no more than k narrowed intervals can have a common point, we have no more than k undecided values x_i . Trying all possible combinations of lower and upper endpoints for these $\leq k$ values requires $\leq 2^k$ steps.

Thus, the overall number of steps is $O(2^k \cdot n^2)$. Since k is a constant, the overall number of steps is thus $O(n^2)$.

The theorem is proven.

Proof of Theorem 5.1

1°. Similarly to the proof of Theorem 3.1, we reduce a subset problem to the problem of computing \bar{V} .

Each instance of the subset problem is as follows: given n positive integers s_1, \dots, s_n , to check whether there exist signs $\eta_i \in \{-1, +1\}$ for which the signed sum $\sum_{i=1}^n \eta_i \cdot s_i$ equals 0.

We will show that this problem can be reduced to the problem of computing \bar{C} , i.e., that to every instance (s_1, \dots, s_n) of the subset problem \mathcal{P}_0 , we can put into correspondence such an instance of the \bar{C} -computing problem that based on its solution, we can easily check whether the desired signs exist.

As this instance, we take the instance corresponding to the intervals $[\underline{x}_i, \bar{x}_i] = [\underline{y}_i, \bar{y}_i] = [-s_i, s_i]$. We want to show that for the corresponding problem, $\bar{C} = C_0$ (where is the same as in the proof of Theorem 3.1) if and only if there exist signs η_i for which $\sum \eta_i \cdot s_i = 0$.

2°. Let us first show that in all cases, $\bar{C} \leq C_0$.

Indeed, it is known that the sample covariance C is bounded by the product $\sigma_x \cdot \sigma_y$ of sample standard deviations $\sigma_x = \sqrt{V_x}$ and $\sigma_y = \sqrt{V_y}$ of x and y . In the proof of Theorem 3.1, we have already proven that the sample variance V_x of the values x_1, \dots, x_n satisfies the inequality $V_x \leq C_0$; similarly, the sample variance V_y of the values y_1, \dots, y_n satisfies the inequality $V_y \leq C_0$. Hence, $C \leq \sigma_x \cdot \sigma_y \leq \sqrt{C_0} \cdot \sqrt{C_0} = C_0$. In other words, every possible value C of the sample covariance is smaller than or equal to C_0 . Thus, the largest of these possible values, i.e., \bar{C} , also cannot exceed C_0 , i.e., $\bar{C} \leq C_0$.

3°. Let us now show that if $\bar{C} = C_0$, then the desired signs exist.

Indeed, if $\bar{C} = C$, this means that for the corresponding values of x_i and y_i , the sample covariance C is equal to C_0 , i.e.,

$$C = C_0 = \frac{1}{n-1} \cdot \sum_{i=1}^n s_i^2.$$

On the other hand, we have shown that in all cases (and in this case in particular), $C \leq \sigma_x \cdot \sigma_y \leq \sqrt{C_0} \cdot \sqrt{C_0} = C_0$. If $\sigma_x < \sqrt{C_0}$, then we would have $C < C_0$. So, if $C = C_0$, we have $\sigma_x = \sigma_y = \sqrt{C_0}$, i.e., $V_x = V_y = C_0$. We have already shown, in the proof of Theorem 3.1, that in this case the desired signs exist.

4°. To complete the proof of Theorem 5.1, we must show that, vice versa, if the desired signs η_i exist, then $\bar{C} = C_0$.

Indeed, in this case, for $x_i = y_i = \eta_i \cdot s_i$, we have $\bar{x} = \bar{y} = 0$ and $x_i \cdot y_i = s_i^2$, hence

$$C = \frac{1}{n-1} \cdot \sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y}) = \frac{1}{n-1} \cdot \sum_{i=1}^n s_i^2 = C_0.$$

The theorem is proven.

Proof of Theorem 5.2

This proof is similar to the proof of Theorem 5.1, with the only difference that in this case, we use the other part of the inequality $|C| \leq \sigma_x \cdot \sigma_y$, namely, that $C \geq -\sigma_x \cdot \sigma_y$, and in the last part of the proof, we take $y_i = -x_i$.

Proof of Theorem 5.3

1°. Similarly to the proof of Theorems 3.1 and 5.1, we reduce a subset problem to the problem of computing \bar{V} .

Each instance of the subset problem is as follows: given m positive integers s_1, \dots, s_m , to check whether there exist signs $\eta_i \in \{-1, +1\}$ for which the signed sum $\sum_{i=1}^m \eta_i \cdot s_i$ equals 0.

We will show that this problem can be reduced to the problem of computing $\bar{\rho}$, i.e., that to every instance (s_1, \dots, s_m) of the subset problem \mathcal{P}_0 , we can put into correspondence such an instance of the $\bar{\rho}$ -computing problem that based on its solution, we can easily check whether the desired signs exist.

As this instance, we take the instance corresponding to the following intervals:

- $n = m + 2$ (note the difference between this reduction and reductions from the proofs of Theorems 3.1 and 5.1, where we have $n = m$);
- $[\underline{x}_i, \bar{x}_i] = [-s_i, s_i]$ and $\mathbf{y}_i = [0, 0]$ for $i = 1, \dots, m$;
- $\mathbf{x}_{m+1} = \mathbf{y}_{m+2} = [1, 1]$; $\mathbf{x}_{m+2} = \mathbf{y}_{m+1} = [-1, -1]$.

Like in the proof of Theorem 3.1, we define C_1 as

$$C_1 = \sum_{i=1}^m s_i^2. \quad (6.27)$$

We will prove that for the corresponding problem, $\bar{\rho} = -\sqrt{\frac{2}{C_1 + 2}}$ if and only if there exist signs η_i for which $\sum \eta_i \cdot s_i = 0$.

2°. The correlation coefficient is defined as $\rho = C / \sqrt{V_x} \cdot \sqrt{V_y}$. To find the range for ρ , it is therefore reasonable to first find ranges for C , V_x , and V_y .

3°. Of these three, the variance V_y is the easiest to compute because there is no interval uncertainty in y_i at all. For y_i , we have $\bar{y} = 0$ and therefore,

$$V_y = \frac{1}{n-1} \cdot \sum_{i=1}^n y_i^2 - \frac{n}{n-1} \cdot (\bar{y})^2 = \frac{2}{n-1} = \frac{2}{m+2}. \quad (6.28)$$

4°. To find the range for the covariance, we will use the known equivalent formula

$$C = \frac{1}{n-1} \cdot \sum_{i=1}^n x_i \cdot y_i - \frac{n}{n-1} \cdot \bar{x} \cdot \bar{y}. \quad (6.29)$$

Since $\bar{y} = 0$, the second sum in this formula is 0, so C is equal to the first sum. In this first sum, the first m terms are 0's because for $i = 1, \dots, m$, we have $y_i = 0$. The only non-zero terms correspond to $i = m+1$ and $i = m+2$, so

$$C = -\frac{2}{n-1} = -\frac{2}{m+2}. \quad (6.30)$$

5°. Substituting the formulas (6.28) and (6.30) into the definition (5.1) of sample correlation, we conclude that

$$\rho = -\frac{\frac{2}{m+1}}{\sqrt{\frac{2}{m+1}} \cdot \sqrt{V_x}} = -\sqrt{\frac{2}{(m+1) \cdot V_x}}. \quad (6.31)$$

Therefore, the sample correlation ρ attains its maximum $\bar{\rho}$ if and only if the sample variance V_x takes the largest possible value \bar{V}_x :

$$\bar{\rho} = -\sqrt{\frac{2}{(m+1) \cdot \bar{V}_x}}. \quad (6.32)$$

Thus, if we can know $\bar{\rho}$, we can reconstruct \bar{V}_x as

$$\bar{V}_x = \frac{2}{(m+1) \cdot (\bar{\rho})^2}. \quad (6.33)$$

In particular, the desired value $\bar{\rho} = -\sqrt{\frac{2}{C_1+2}}$ corresponds to $\bar{V}_x = \frac{C_1+2}{m+1}$. Therefore, to complete our proof, we must show that $\bar{V}_x = \frac{C_1+2}{m+1}$ if and only if there exist signs η_i for which $\sum \eta_i \cdot s_i = 0$.

6°. Similarly to the proof of Theorem 3.1, we will use the equivalent expression (6.6) for the sample variance V_x ; we will slightly reformulate this expression by substituting the definition of \bar{x} into it:

$$V_x = \frac{1}{n-1} \cdot \sum_{i=1}^n x_i^2 - \frac{1}{n \cdot (n-1)} \cdot \left(\sum_{i=1}^n x_i \right)^2. \quad (6.34)$$

We can (somewhat) simplify this expression by substituting the values $n = m+2$, $x_{m+1} = 1$, and $x_{m+2} = -1$. We have

$$\sum_{i=1}^n x_i = \sum_{i=1}^m x_i + x_{m+1} + x_{m+2} = \sum_{i=1}^m x_i$$

and

$$\sum_{i=1}^n x_i^2 = \sum_{i=1}^m x_i^2 + x_{m+1}^2 + x_{m+2}^2 = \sum_{i=1}^m x_i + 2.$$

Therefore,

$$V_x = \frac{1}{m+1} \cdot \sum_{i=1}^m x_i^2 + \frac{2}{m+1} - \frac{1}{(m+1) \cdot (m+2)} \cdot \left(\sum_{i=1}^m x_i \right)^2. \quad (6.35)$$

Similarly to the proof of Theorem 3.1, we can show that always $V_x \leq \frac{C_1+2}{m+1}$, and that $\bar{V}_x = \frac{C_1+2}{m+1}$ if and only if there exist the signs η_i for which $\sum \eta_i \cdot s_i = 0$.

The theorem is proven.

Proof of Theorem 5.4

This proof is similar to the proof of Theorem 5.3, with the only difference that we take $y_{m+1} = 1$ and $y_{m+2} = -1$. In this case,

$$C = \frac{2}{m+2},$$

hence

$$\rho = \sqrt{\frac{2}{(m+1) \cdot V_x}},$$

and so the largest possible value of V_x corresponds to the smallest possible value of ρ .

Acknowledgments

This work was supported in part by NASA under cooperative agreement NCC5-209 and grant NCC 2-1232, by NSF grants CDA-9522207, ERA-0112968, and 9710940 Mexico/Conacyt, by the Future Aerospace Science and Technology Program (FAST) Center for Structural Integrity of Aerospace Systems, effort sponsored by the Air Force Office of Scientific Research, Air Force Materiel Command, USAF, under grants numbers F49620-95-1-0518 and F49620-00-1-0365, by grant No. W-00016 from the U.S.-Czech Science and Technology Joint Fund, by IEEE/ACM SC2001 Minority Serving Institutions Participation Grant, and by Small Business Innovation Research grant 9R44CA81741 to Applied Biomathematics from the National Cancer Institute (NCI), a component of the National Institutes of Health (NIH). The opinions expressed herein are those of the author(s) and not necessarily those of NASA, NSF, AFOSR, NCI, or the NIH.

The authors are greatly thankful to Daniel E. Cooke, Michael Gelfond, Renata Maria C. R. de Souza, and to the anonymous referees for very useful suggestions.

References

- [1] Th. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to algorithms*, MIT Press, Cambridge, MA, and Mc-Graw Hill Co., N.Y., 2001.
- [2] S. Ferson, W. Root and R. Kuhn, *RAMAS Risk Calc: Risk Assessment with Uncertain Numbers*, Applied Biomathematics, Setauket, NY, 1999.
- [3] W. A. Fuller, *Measurement error models*, J. Wiley & Sons, New York, 1987.
- [4] M. R. Garey and S. D. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*, W.H. Freeman and Company, New York, 1979.
- [5] V. Kreinovich, A. Lakeyev, J. Rohn, and P. Kahl, *Computational complexity and feasibility of data processing and interval computations*, Kluwer, Dordrecht, 1997.
- [6] V. P. Kuznetsov, *Interval statistical models*, Moscow, Radio i Svyaz Publ., 1991 (in Russian).
- [7] R. Osegueda, V. Kreinovich, L. Potluri, and R. Aló, “Non-Destructive Testing of Aerospace Structures: Granularity and Data Mining Approach”, *Proceedings of FUZZ-IEEE'2002*, Honolulu, Hawaii, May 12-17, 2002 (to appear).
- [8] C. H. Papadimitriou, *Computational Complexity*, Addison Wesley, San Diego, 1994.
- [9] S. Rabinovich, *Measurement Errors: Theory and Practice*, American Institute of Physics, New York, 1993.
- [10] S. A. Vavasis, *Nonlinear optimization: complexity issues*, Oxford University Press, N.Y., 1991.
- [11] H. M. Wadsworth, Jr. (eds.), *Handbook of statistical methods for engineers and scientists*, McGraw-Hill Publishing Co., New York, 1990.
- [12] P. Walley, *Statistical reasoning with imprecise probabilities*, Chapman and Hall, N.Y., 1991.
- [13] G. W. Walster, “Philosophy and practicalities of interval arithmetic”, In: R. E. Moore (ed.), *Reliability in Computing*, 1988, pp. 307–323.