

Computing Population Variance and Entropy under Interval Uncertainty: Linear-Time Algorithms

Gang Xiang, Martine Ceberio, and Vladik Kreinovich
Department of Computer Science
University of Texas at El Paso, El Paso, TX 79968, USA,
gxiang@utep.edu, mceberio@cs.utep.edu, vladik@utep.edu

Abstract

In statistical analysis of measurement results, it is often necessary to compute the range $[\underline{V}, \overline{V}]$ of the population variance $V = \frac{1}{n} \cdot \sum_{i=1}^n (x_i - E)^2$ (where $E = \frac{1}{n} \cdot \sum_{i=1}^n x_i$) when we only know the intervals $[\tilde{x}_i - \Delta_i, \tilde{x}_i + \Delta_i]$ of possible values of the x_i . While \underline{V} can be computed efficiently, the problem of computing \overline{V} is, in general, NP-hard. In our previous paper “Population Variance under Interval Uncertainty: A New Algorithm” (*Reliable Computing*, 2006, Vol. 12, No. 4, pp. 273–280), we showed that in a practically important case, we can use constraints techniques to compute \overline{V} in time $O(n \cdot \log(n))$. In this paper, we provide new algorithms that compute \underline{V} and, for the above case, \overline{V} in linear time $O(n)$.

Similar linear-time algorithms are described for computing the range of the entropy $S = -\sum_{i=1}^n p_i \cdot \log(p_i)$ when we only know the intervals $\mathbf{p}_i = [\underline{p}_i, \overline{p}_i]$ of possible values of probabilities p_i .

1 Computing Population Variance under Interval Uncertainty: Formulation of the Problem

Statistical analysis is important. Once we have n measurement results x_1, \dots, x_n , the traditional statistical analysis starts with computing the standard statistics such as population mean $E = \frac{1}{n} \cdot \sum_{i=1}^n x_i$ and

population variance $V = M - E^2$, where $M \stackrel{\text{def}}{=} \frac{1}{n} \cdot \sum_{i=1}^n x_i^2$; see, e.g., [14].

These values are useful, e.g., in detecting *outliers*: once we know the mean E and the standard deviation $\sigma \stackrel{\text{def}}{=} \sqrt{V}$ of the normal values, we can determine outliers as values x_i for which $|x_i - E| \gg \sigma$, i.e., $|x_i - E| \geq k_0 \cdot \sigma$ for some k_0 (usually, $k_0 = 2, 3$, or 6).

Outliers are important in many application areas: in non-destructive testing, outliers indicate possible faults; in geophysics, outliers should be identified as possible locations of minerals; in medicine, outliers indicate possible illnesses, etc.

The more data point we take, the more accurate the resulting estimates for E and V . Thus, in many practical applications, we process large amount of data: in geophysics, we process thousands and millions data points; in processing census data, we process data about millions of people, etc.

Interval uncertainty. In many real-life situations, due to measurement uncertainty, instead of the actual values x_i of the measured quantity, we only have intervals $\mathbf{x}_i = [\underline{x}_i, \overline{x}_i]$ of possible values of x_i [7, 14]. Usually, the interval \mathbf{x}_i has the form $[\tilde{x}_i - \Delta_i, \tilde{x}_i + \Delta_i]$, where \tilde{x}_i is the measurement result, and Δ_i is the known upper bound on the absolute value $|\Delta x_i|$ of the (unknown) measurement error $\Delta x_i \stackrel{\text{def}}{=} \tilde{x}_i - x_i$: $|\Delta x_i| \leq \Delta_i$.

Another source of interval uncertainty is the existence of *detection limits* for different sensors: if a sensor, e.g., did not detect any ozone, this means that the ozone concentration is below its detection limit DL , i.e., in the interval $[0, DL]$.

One more source of interval uncertainty is *discretization*: to study the effect of a pollutant on the fish, we check on the fish daily; if a fish was alive on Day 5 but dead on Day 6, then the only information about the lifetime of this fish is that it is somewhere within the interval $[5, 6]$; we have no information about the distribution of different values in this interval.

Yet another source of interval uncertainty is *privacy*. In biomedical systems, statistical analysis of the data often leads to improvements in medical recommendations; however, to maintain privacy, we do not want to use the exact values of the patient's parameters. Instead, for each parameter, we select fixed values (thresholds), and for each patient, we only keep the corresponding range. For example, instead of keeping the exact age, we only record whether the age is between 0 and 10, 10 and 20, 20 and 30, etc.

Finally, intervals occur if instead of measurements, we use *expert estimates* for difficult-to-measure quantities: experts can rarely describe exact values of the physical quantities; at best, they can provide the bounds on the possible values, i.e., intervals which contain the (unknown) actual value of the quantity of interest.

In statistical analysis, it is necessary to take into account interval uncertainty. As we have mentioned, in many practical situations, it is desirable to know the mean E and the variance V . In case of interval uncertainty, different values $x_i \in \mathbf{x}_i$ lead, in general, to different values of E and V . It is therefore desirable to compute the ranges $\mathbf{E} = [\underline{E}, \overline{E}]$ and $\mathbf{V} = [\underline{V}, \overline{V}]$ of possible values of E and V when $x_i \in \mathbf{x}_i$.

Example of practical applications: in brief. Interval ranges for statistical characteristics have been successfully used in geophysics [12, 13], in environmental science, and in many other application areas; see [10, 11] and references therein.

Computing the range of variance under interval uncertainty: what is known. Since the population mean E is a monotonic function of its n variables x_1, \dots, x_n , its range can be easily computed as

$$\mathbf{E} = \left[\frac{1}{n} \cdot \sum_{i=1}^n \underline{x}_i, \frac{1}{n} \cdot \sum_{i=1}^n \overline{x}_i \right].$$

The population variance $V(x_1, \dots, x_n)$ is, in general, not a monotonic function of its variables x_i ; so, we cannot easily indicate the values at which this function attains its minimum and maximum. It can be easily checked that population variance is a convex function; thus, we can use known efficient (polynomial time) algorithms of convex optimization to find the minimum \underline{V} of this convex function $V(x_1, \dots, x_n)$ on the convex box $\mathbf{x}_1 \times \dots \times \mathbf{x}_n$. For this particular convex function, it is possible to describe algorithms which are faster than in the general convex case: namely, we can compute the lower bound \underline{V} in time $O(n \cdot \log(n))$; see, e.g., [11].

For the upper bound \overline{V} , the situation is more complicated. Specifically, it is known that the maximum of a convex function on a convex set is attained at one of its extreme points. Thus, the maximum \overline{V} is attained at one of the extreme points of the convex box $\mathbf{x}_1 \times \dots \times \mathbf{x}_n$, i.e., when for every i , the variable x_i is equal to one of the endpoints: $x_i = \underline{x}_i$ or $x_i = \overline{x}_i$. In general, there are 2^n combinations of n such endpoints; so, we can compute \overline{V} by finding the maximum of 2^n corresponding values. The computation time for this computation grows exponentially with the size n of the problem.

It is known that in general, computing \overline{V} is an NP-hard problem [5, 11]. This NP-hardness result means that (unless $P=NP$) no general algorithm is possible that would always compute \overline{V} in feasible (polynomial) time – i.e., crudely speaking, that in the worst case, the exponential computation time is inevitable.

It is, however, possible to compute \overline{V} in polynomial time in some practically reasonable cases.

In general, we can have both high-accuracy data points (e.g., points with narrow uncertainty intervals) and low-accuracy data points (e.g., points with much wider uncertainty intervals). For example, in addition to accurate measurements that provide narrow intervals for the values of the desired quantity, we may have expert estimates that come from expert estimates. It is well known that in statistics, if we have a large number of high-accuracy data points, then the additional information provided by low-accuracy data

points is negligible; as a result, usually, in statistical analysis, only high-accuracy data points are taken into account. With this practice in mind, it is reasonable to restrict ourselves to the case when all the intervals are approximately of the same width. This happens, e.g., when all measurements have been done by a single measuring instrument (or by several measuring instruments of the same type).

How can we describe this mathematically? A clear indication that we have two measuring instruments (MI) of different quality is that one interval is a proper subset of the other one: $[\underline{x}_i, \bar{x}_i] \subseteq (\underline{x}_j, \bar{x}_j)$. Thus, the condition that all data points are of the same accuracy means can be formalized as a requirement that no interval is a proper subinterval of another one.

This property holds for other sources of interval uncertainty: e.g., for detection limits, we have intervals of the type $[0, DL_i]$ for which $[0, DL_i] \not\subseteq (0, DL_j)$; in the privacy case, we have intervals $[b_k, b_{k+1}]$ between two consecutive thresholds, none of which is a proper subset of the other, etc.

In this practically useful case, there exists an algorithm for computing \bar{V} in time $O(n \cdot \log(n))$ [3]. Actually, this algorithm works in a more general case, when the above “proper subset” property holds only for the “narrowed” intervals: namely, for the case when no “narrowed interval” $[x_i^-, x_i^+]$, where $x^- \stackrel{\text{def}}{=} \tilde{x}_i - \frac{\Delta_i}{n}$ and $x_i^+ \stackrel{\text{def}}{=} \tilde{x}_i + \frac{\Delta_i}{n}$, is a proper subinterval of the interior of another narrowed interval, i.e., when $|\tilde{x}_i - \tilde{x}_j| \geq \frac{|\Delta_i - \Delta_j|}{n}$ for all $i \neq j$.

In this paper, we describe two new linear-time algorithms:

- a linear-time algorithm that computes \underline{V} for all possible intervals, and
- a linear-time algorithm that computes \bar{V} for intervals that satisfy the above subset property.

2 Linear-Time Algorithm for Computing \bar{V}

Our new algorithm is based on the known fact that we can compute the median of a set of n elements in linear time (see, e.g., [2]); our use of median is similar to the one from [1, 6]. Let us first describe the algorithm itself; in the Appendix, we provide a justification for this algorithm.

For simplicity, let us first consider the case when all the intervals are non-degenerate, i.e., when $\Delta_i > 0$ for all i .

The proposed algorithm is iterative. At each iteration of this algorithm, we have three sets:

- the set I^- of all the indices i from 1 to n for which we already know that for the optimal vector x , we have $x_i = \underline{x}_i$;
- the set I^+ of all the indices j for which we already know that for the optimal vector x , we have $x_j = \bar{x}_j$;
- the set $I = \{1, \dots, n\} - I^- - I^+$ of the indices i for which we are still undecided.

In the beginning, $I^- = I^+ = \emptyset$ and $I = \{1, \dots, n\}$. At each iteration, we also update the values of two auxiliary quantities $E^- \stackrel{\text{def}}{=} \sum_{i \in I^-} \underline{x}_i$ and $E^+ \stackrel{\text{def}}{=} \sum_{j \in I^+} \bar{x}_j$. In principle, we could compute these values by computing these sums, but to speed up computations, on each iteration, we update these two auxiliary values in a way that is faster than re-computing the corresponding two sums. Initially, since $I^- = I^+ = \emptyset$, we take $E^- = E^+ = 0$.

At each iteration, we do the following:

- first, we compute the median m of the set I (median in terms of sorting by \tilde{x}_i);
- then, by analyzing the elements of the undecided set I one by one, we divide them into two subsets $P^- = \{i : \tilde{x}_i \leq \tilde{x}_m\}$ and $P^+ = \{j : \tilde{x}_j > \tilde{x}_m\}$;
- we compute $e^- = E^- + \sum_{i \in P^-} \underline{x}_i$ and $e^+ = E^+ + \sum_{j \in P^+} \bar{x}_j$;
- if $n \cdot x_m^- < e^- + e^+$, then we replace I^- with $I^- \cup P^-$, E^- with e^- , and I with P^+ ;

- if $n \cdot x_m^- > e^- + e^+$, then we replace I^+ with $I^+ \cup P^+$, E^+ with e^+ , and I with P^- ;
- if $n \cdot x_m^- = e^- + e^+$, then we replace I^- with $I^- \cup P^-$, I^+ with $I^+ \cup P^+$, and I with \emptyset .

At each iteration, the set of undecided indices is divided in half. Iterations continue until all indices are decided, after which we return, as \bar{V} , the value of the population variance for the vector x for which $x_i = \underline{x}_i$ for $i \in I^-$ and $x_j = \bar{x}_j$ for $j \in I^+$.

Comments.

- This same algorithm can be easily applied if one of the intervals consists of a single point only: this value is plugged in and the variable is eliminated.
- For readers' convenience, all the proofs – that the algorithms are correct and that they require linear time – are placed in a special Appendix.
- As with all asymptotic results, a natural question arises: how practical is the new linear time $O(n)$ algorithm? For which n is it better than the known $O(n \cdot \log(n))$ algorithm for computing \bar{V} ? In general, the answer depends on the constants in the corresponding asymptotics. The constant for the known $O(n \cdot \log(n))$ algorithm is ≈ 1 . As one can see from the proof, for our new algorithm, the constant is the same as for known linear time algorithm for computing the median, i.e., it is ≈ 20 [2]; thus, the new algorithm is better when $\log_2(n) > 20$, i.e., when $n > 10^6$. We have mentioned that in many practical applications we do need to process millions of data points; in such applications, the new algorithm for computing \bar{V} is indeed faster.

3 Linear-Time Algorithm for Computing \bar{V}

The proposed algorithm is iterative. At each iteration of this algorithm, we have three sets:

- the set J^- of all the endpoints \underline{x}_i and \bar{x}_j for which we already know that for the optimal vector x , we have, correspondingly, $x_i \neq \underline{x}_i$ (for \underline{x}_i) or $x_j = \bar{x}_j$ (for \bar{x}_j);
- the set J^+ of all the endpoints \underline{x}_i and \bar{x}_j for which we already know that for the optimal vector x , we have, correspondingly, $x_i = \underline{x}_i$ (for \underline{x}_i) or $x_j \neq \bar{x}_j$ (for \bar{x}_j);
- the set J of the endpoints \underline{x}_i and \bar{x}_j for which we have not yet decided whether these endpoints appear in the optimal vector x .

In the beginning, $J^- = J^+ = \emptyset$ and J is the set of all $2n$ endpoints. At each iteration, we also update the values $N^- = \#(J^-)$, $N^+ = \#(J^+)$, $E^- = \sum_{\bar{x}_j \in J^-} \bar{x}_j$, and $E^+ = \sum_{\underline{x}_i \in J^+} \underline{x}_i$. Initially, $N^- = N^+ = E^- = E^+ = 0$.

At each iteration, we do the following:

- first, we compute the median m of the set J ;
- then, by analyzing the elements of the undecided set J one by one, we divide them into two subsets

$$Q^- = \{x \in J : x \leq m\}; \quad Q^+ = \{x \in J : x > m\};$$

we also compute $m^+ = \min\{x : x \in Q^+\}$;

- we compute $e^- = E^- + \sum_{\bar{x}_j \in Q^-} \bar{x}_j$, $e^+ = E^+ + \sum_{\underline{x}_i \in Q^+} \underline{x}_i$,

$$n^- = N^- + \#\{\bar{x}_j \in Q^-\}, \quad n^+ = N^+ + \#\{\underline{x}_i \in Q^+\},$$

$$\text{and } r = \frac{e^- + e^+}{n^- + n^+};$$

- if $r < m$, then we replace J^- with $J^- \cup Q^-$, E^- with e^- , J with Q^+ , and N^- with n^- ;
- if $r > m^+$, then we replace J^+ with $J^+ \cup Q^+$, E^+ with e^+ , J with P^- , and N^+ with n^+ ;
- if $m \leq r \leq m^+$, then we replace J^- with $J^- \cup Q^-$, J^+ with $J^+ \cup Q^+$, J with \emptyset , E^- with e^- , E^+ with e^+ , N^- with n^- , and N^+ with n^+ .

At each iteration, the set of undecided indices is divided in half. Iterations continue until all indices are decided, after which we return, as \underline{V} , the value of the population variance for the vector x for which:

- $x_j = \bar{x}_j$ for indices j for which $\bar{x}_j \in J^-$,
- $x_i = \underline{x}_i$ for indices i for which $\underline{x}_i \in J^+$, and
- $x_i = r$ for all other indices i .

4 Computing Entropy under Interval Uncertainty: Formulation of the Problem

For a probability distribution with probabilities p_1, \dots, p_n , $\sum p_i = 1$, the amount of uncertainty can be described by Shannon's entropy $S = -\sum_{i=1}^n p_i \cdot \log(p_i)$. In practice, we sometimes only know the intervals $\mathbf{p}_i = [\underline{p}_i, \bar{p}_i]$ of possible values of p_i . Different values $p_i \in \mathbf{p}_i$ lead to different S ; it is therefore desirable to find the range $[\underline{S}, \bar{S}]$ of the entropy S [8].

Since the function S is concave, computation of \bar{S} is feasible [9, 15]; however, computing \underline{S} is NP-hard [16]. For the case when no interval $[\underline{p}_i, \bar{p}_i]$ is a proper subset of the interior of another interval \mathbf{p}_j , we have proposed, in [16], an $O(n \cdot \log(n))$ algorithm for computing \underline{S} .

In this paper, we describe two new linear-time algorithms:

- a linear-time algorithm that computes \bar{S} for all possible intervals, and
- a linear-time algorithm that computes \underline{S} for intervals that satisfy the above subset property.

Comment. The new algorithms are similar to the above algorithms for computing \underline{V} and \bar{V} . The main difference between the algorithms for variance and the algorithms for entropy is that variance $V(x_1, \dots, x_n)$ is defined for all possible combinations $x_i \in [\underline{x}_i, \bar{x}_i]$, but the entropy is only defined for values $p_i \in [\underline{p}_i, \bar{p}_i]$ that satisfy the additional constraint $\sum_{i=1}^n p_i = 1$. As a result, while the maximum \bar{V} of the variance V is attained when *each* value x_i attains one of its endpoints, the minimum \underline{S} of the entropy S is attained, in general, when *all but one* values are endpoints: it may not be possible to have $\sum_{i=1}^n p_i = 1$ if we only use endpoints.

5 Linear-Time Algorithm for Computing \underline{S}

The proposed algorithm for computing \underline{S} is similar to the algorithm for computing \bar{V} ; the only difference is in the replacement part: once we computed m , P^- , P^+ , e^- , and e^+ , we do the following:

- if $e^- + e^+ > 1$, then we replace I^- with $I^- \cup P^-$, E^- with e^- , and I with P^+ ;
- if $e^- + e^+ + 2\Delta_m < 1$, then we replace I^+ with $I^+ \cup P^+$, E^+ with e^+ , and I with P^- ;
- finally, if $e^- + e^+ \leq 1 \leq e^- + e^+ + 2\Delta_m$, then we replace I^- with $I^- \cup (P^- - \{m\})$, I^+ with $I^+ \cup P^+$, I with $\{m\}$, E^- with $e^- - \underline{p}_m$, and E^+ with e^+ .

When computing \bar{V} , iterations continued until $I = \emptyset$. For \underline{S} , iterations continue until we have only one undecided index $I = \{k\}$, after which we return, as \underline{S} , the value of the entropy for the vector p for which $p_i = \underline{p}_i$ for $i \in I^-$, $x_j = \bar{p}_j$ for $j \in I^+$, and $p_k = 1 - e^- - e^+$ for the remaining value k .

Comment: a similar linear-time algorithm can be used to compute the expected value under interval uncertainty. In addition to computing the range for entropy, it is often useful to compute the range $[\underline{a}, \bar{a}]$ of the expected value $a = \sum a_i \cdot p_i$ of a known variable (a_1, \dots, a_n) under the constraints $p_i \in [\underline{p}_i, \bar{p}_i]$ and $\sum_{i=1}^n p_i = 1$. For this problem, linear-time algorithms are known; see, e.g., [1, 6]. Let us show that this problem can be also solved by a simple modification of the above algorithm.

It is known that the smallest possible value \underline{a} of the linear form $\sum_{i=1}^n a_i \cdot p_i$ under given constraints is equal to $-\bar{b}$, where \bar{b} is the largest possible value of the form $\sum_{i=1}^n b_i \cdot p_i$, with $b_i = -a_i$. Thus, it is sufficient to describe how to compute \bar{a} .

For that, we follow the above iterative algorithm while it computes I^- and I^+ . We continue iterations until we have only one undecided index $I = \{k\}$, after which we return, as \bar{a} , the value of the linear function $\sum_{i=1}^n a_i \cdot p_i$ for the vector p for which $p_i = \underline{p}_i$ for $i \in I^-$, $x_j = \bar{p}_j$ for $j \in I^+$, and $p_k = 1 - e^- - e^+$ for the remaining value k .

6 Linear-Time Algorithm for Computing \bar{S}

The proposed algorithm is similar to the algorithm for computing \underline{V} ; the only difference is that for computing \bar{S} , at each iteration, instead of computing $r = \frac{e^- + e^+}{n^- + n^+}$, we compute $r = \frac{1 - e^- - e^+}{1 - n^- - n^+}$.

7 Open Questions

Can similar linear-time algorithms be proposed for computing the endpoints of the intervals for the quantities $E - \alpha \cdot \sqrt{V}$ and $E + \alpha \cdot \sqrt{V}$ – which are important in detecting outliers [4, 10]? for computing other statistical characteristics – like moments or covariance?

Acknowledgments

This work was supported in part by NASA under cooperative agreement NCC5-209, NSF grants EAR-0225670 and DMS-0532645, Star Award from the University of Texas System, and Texas Department of Transportation grant No. 0-5453.

The authors are greatly thankful to Pim van den Broek for important discussions and references, to Scott Ferson for valuable discussions, and to the anonymous referees for important suggestions.

References

- [1] P. van der Broek and J. Noppen, “Fuzzy weighted average: alternative approach”, *Proceedings of the 25th International Conference of the North American Fuzzy Information Processing Society NAFIPS’2006*, Montreal, Quebec, Canada, June 3–6, 2006.
- [2] Th. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms*, MIT Press, Cambridge, MA, 2001.
- [3] E. Dantsin, V. Kreinovich, A. Wolpert, and G. Xiang, “Population Variance under Interval Uncertainty: A New Algorithm”, *Reliable Computing*, 2006, Vol. 12, No. 4, pp. 273–280.
- [4] E. Dantsin, A. Wolpert, M. Ceberio, G. Xiang, and V. Kreinovich, “Detecting Outliers under Interval Uncertainty: A New Algorithm Based on Constraint Satisfaction”, *Proceedings of the International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems IPMU’06*, Paris, France, July 2–7, 2006, pp. 802–809.

- [5] S. Ferson, L. Ginzburg, V. Kreinovich, L. Longpré, and M. Aviles, Exact Bounds on Finite Populations of Interval Data, *Reliable Computing*, 2005, Vol. 11, No. 3, pp. 207–233.
- [6] P. Hansen, M. V. P. de Aragao, and C. C. Ribeiro, “Hyperbolic 0-1 programming and optimization in information retrieval”, *Math. Programming*, 1991, Vol. 52, pp. 255–263.
- [7] L. Jaulin, M. Kieffer, O. Didrit, and E. Walter, *Applied interval analysis: with examples in parameter and state estimation, robust control and robotics*, Springer Verlag, London, 2001.
- [8] G. J. Klir, *Uncertainty and Information: Foundations of Generalized Information Theory*, J. Wiley, Hoboken, New Jersey, 2005.
- [9] V. Kreinovich, “Maximum entropy and interval computations”, *Reliable Computing*, 1996, Vol. 2, No. 1, pp. 63–79.
- [10] V. Kreinovich, L. Longpré, P. Patangay, S. Ferson, and L. Ginzburg, “Outlier Detection Under Interval Uncertainty: Algorithmic Solvability and Computational Complexity”, *Reliable Computing*, 2005, Vol. 11, No. 1, pp. 59–76.
- [11] V. Kreinovich, G. Xiang, S. A. Starks, L. Longpré, M. Ceberio, R. Araiza, J. Beck, R. Kandathi, A. Nayak, R. Torres, and J. Hajagos, “Towards combining probabilistic and interval uncertainty in engineering calculations: algorithms for computing statistics under interval uncertainty, and their computational complexity”, *Reliable Computing*, 2006, Vol. 12, No. 6, pp. 471–501.
- [12] P. Nivlet, F. Fournier, and J. Royer, 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, 2001, 1644–1647.
- [13] P. Nivlet, F. Fournier, and J. Royer, 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, 2001, paper SPE-71327.
- [14] S. Rabinovich, *Measurement Errors: Theory and Practice*, American Institute of Physics, New York, 1993.
- [15] S. A. Vavasis, *Nonlinear Optimization: Complexity Issues*, Oxford Science, New York, 1991.
- [16] G. Xiang, O. Kosheleva, and G. J. Klir, “Estimating information amount under interval uncertainty: algorithmic solvability and computational complexity”, *Proceedings of the International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems IPMU’06*, Paris, France, July 2–7, 2006, pp. 840–847.

Appendix: Proofs

Proof that the new algorithm for computing \bar{V} requires linear time. At each iteration, computing median requires linear time, and all other operations with I require time t linear in the number of elements $|I|$ of I : $t \leq C \cdot |I|$ for some C . We start with the set I of size n ; on the next iteration, we have a set of size $n/2$, then $n/4$, etc. Thus, the overall computation time is $\leq C \cdot (n + n/2 + n/4 + \dots) \leq C \cdot 2n$, i.e., linear in n .

Proof that under the subset property, the new algorithm always computes \bar{V} . Similarly to [11], one can easily show that since no two narrowed intervals are proper subsets of one another, they can be linearly ordered in lexicographic order. In this order, we have $x_1^- \leq x_2^- \leq \dots \leq x_n^-$, $x_1^+ \leq x_2^+ \leq \dots \leq x_n^+$, and, thus, the averages $\tilde{x}_i = (x_i^- + x_i^+)/2$ are also sorted: $\tilde{x}_1 \leq \tilde{x}_2 \leq \dots \leq \tilde{x}_n$.

In [3], we have shown that in this sorting, the value \bar{V} is attained at one of the vectors $x^{(k)} = (\underline{x}_1, \dots, \underline{x}_k, \bar{x}_{k+1}, \dots, \bar{x}_n)$, i.e., that $V = V(x^{(k)})$ for some k .

In [3], we also analyzed the change in $V(x^{(k)})$ when we replace $x^{(k)}$ with $x^{(k-1)}$, i.e., when we replace \underline{x}_k with $\bar{x}_k = \underline{x}_k + 2\Delta_k$; we have shown that $V_{k-1} - V_k = \frac{4\Delta_k}{n} \cdot (x_k^- - E_k)$, where $E_k \stackrel{\text{def}}{=} E(x^{(k)})$.

Hence, $V_{k-1} < V_k$ if and only if $x_k^- < E_k$. Multiplying both sides of this inequality by n , we get an equivalent inequality $x_k^- < n \cdot E_k$, where $n \cdot E_k = \sum_{i=1}^k \underline{x}_i + \sum_{j=k+1}^n \bar{x}_j$. Similarly, $V_{k-1} > V_k$ if and only if $x_k^- > E_k$, and $V_{k-1} = V_k$ if and only if $x_k^- = E_k$.

When we go from k to $k+1$, we replace the larger value \bar{x}_{k+1} in the sum $n \cdot E_k$ by a smaller value \underline{x}_k . Thus, the sequence $n \cdot E_k$ is strictly decreasing with k , while x_k^- is (maybe non-strictly) increasing with k . So, once we have $n \cdot x_k^- < E_k$, i.e., $V_{k-1} < V_k$, these inequalities will hold for smaller k as well. Similarly, once we have $n \cdot x_k^- > E_k$, i.e., $V_{k-1} > V_k$, these inequalities will hold for larger k as well.

Once we have $n \cdot x_k^- = E_k$, i.e., $V_{k-1} = V_k$, then we will have $V_k > V_{k+1} > \dots$ and $V_k = V_{k-1} > V_{k-2} > \dots$, i.e., $V_k = V_{k-1}$ will be the largest value of V .

In other words, the sequence V_k first increases ($V_k > V_{k-1}$ for $k = 1, 2, \dots$) and then starts decreasing ($V_k < V_{k-1}$ for larger k), with one or two top values.

For each m , if $V_{m-1} < V_m$ (i.e., if $n \cdot x_m^- < E_m$), this means that the value k_{\max} corresponding to the maximum of V is $\leq m$; hence for all the indices $\leq m$, we already know that in the optimal vector x , $x_i = \underline{x}_i$. Thus, these indices can be added to the set I^- .

If $V_m > V_{m-1}$ (i.e., if $n \cdot x_m^- > E_m$), this means that the value k_{\max} corresponding to the maximum of V is $> m$; hence for all the indices $> m$, we already know that in the optimal vector x , $x_i = \bar{x}_i$. Thus, these indices can be added to the set I^+ .

Finally, if $V_m = V_{m-1}$ (i.e., if $n \cdot x_m^- = E_m$), then this m is where the maximum is attained.

The algorithm has been justified.

Proof that the new algorithm for computing \underline{V} requires linear time. At each iteration, computing median requires linear time, and all other operations with J require time t linear in the number of elements $|J|$ of J . We start with the set J of size n ; on the next iteration, we have a set of size $n/2$, then $n/4$, etc. Thus, the overall computation time is $\leq C \cdot (n + n/2 + n/4 + \dots) \leq C \cdot 2n$, i.e., linear in n .

Proof that the new algorithm always computes \underline{V} . In [5], we proved that if we sort all $2n$ endpoints into a sequence $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(2n)}$, then for some $k = k_{\min}$, the minimum \underline{V} is attained for the vector x for which:

- for all indices j for which $\bar{x}_j \leq x_{(k)}$, we have $x_j = \bar{x}_j$;
- for all indices i for which $\underline{x}_i \geq x_{(k+1)}$, we have $x_i = \underline{x}_i$;
- for all other indices, we have $x_i = r_k \stackrel{\text{def}}{=} \frac{E_k}{N_k}$, where

$$E_k = \sum_{j: \bar{x}_j \leq x_{(k)}} \bar{x}_j + \sum_{i: \underline{x}_i \geq x_{(k+1)}} \underline{x}_i; \quad N_k = \#\{j : \bar{x}_j \leq x_{(k)}\} + \#\{i : \underline{x}_i \geq x_{(k+1)}\}.$$

It has also been proven that for the optimal k , we have $r_k \in [x_{(k)}, x_{(k+1)}]$.

In general, the condition $x_{(k)} \leq r_k = \frac{E_k}{N_k}$ is equivalent to

$$N_k \cdot x_{(k)} \leq E_k = \sum_{j: \bar{x}_j \leq x_{(k)}} \bar{x}_j + \sum_{i: \underline{x}_i \geq x_{(k+1)}} \underline{x}_i.$$

Subtracting $x_{(k)}$ from each of N_k terms in the right-hand side (RHS), and moving the sum of the resulting non-positive differences into the left-hand side (LHS), we conclude that

$$\sum_{j: \bar{x}_j \leq x_{(k)}} (x_{(k)} - \bar{x}_j) \leq \sum_{i: \underline{x}_i \geq x_{(k+1)}} (\underline{x}_i - x_{(k)}). \quad (1)$$

When we increase k , we get, in general, more terms in the LHS and fewer in the RHS, so LHS (non-strictly) increases, while the RHS non-strictly decreases. So, if the inequality (1) holds for some k , it holds for all smaller values of k as well. Thus, this inequality holds for all k until a certain value k_0 .

Similarly, the condition $x_{(k+1)} \geq r_k = \frac{E_k}{N_k}$ is equivalent to

$$N_k \cdot r_{k+1} \geq \sum_{j:\bar{x}_j \leq x_{(k)}} \bar{x}_j + \sum_{i:\underline{x}_i \geq x_{(k+1)}} \underline{x}_i.$$

Subtracting $x_{(k+1)}$ from each of N_k terms in RHS, and moving the sum of the resulting non-positive differences into LHS, we conclude that

$$\sum_{j:\bar{x}_j \leq x_{(k)}} (x_{(k+1)} - \bar{x}_j) \geq \sum_{i:\underline{x}_i \geq x_{(k+1)}} (\underline{x}_i - x_{(k+1)}). \quad (2)$$

When we increase k , the LHS (non-strictly) increases, while the RHS non-strictly decreases. So, if the inequality (2) holds for some k , it holds for all larger values of k as well. Thus, this inequality holds for all k after a certain value l_0 .

So, both conditions (1) and (2) are satisfied (which is equivalent to the condition $r_k \in [x_{(k)}, x_{(k+1)}]$) either for a single value k_{\min} , or for several sequential values $l_0, l_0 + 1, \dots, k_0$. Let us show that if this condition is satisfied for several sequential values, this simply means that the same minimum \underline{V} is attained for all these values. For that, it is sufficient to show that if both conditions (1) and (2) holds for k and for $k + 1$, then the variance V has the same value for both k and $k + 1$. Indeed, since (1) is true for $k + 1$, we have

$$\sum_{j:\bar{x}_j \leq x_{(k+1)}} (x_{(k+1)} - \bar{x}_j) \leq \sum_{i:\underline{x}_i \geq x_{(k+2)}} (\underline{x}_i - x_{(k+1)}).$$

The LHS of this new inequality is smaller than or equal to the LHS of the inequality (2), and its RHS is larger than or equal to the RHS of the inequality (2). Thus, the only way for both inequalities to hold is when both sides are equal, i.e., when replacing $x_{(k)}$ with $x_{(k+1)}$ and replacing $x_{(k+1)}$ with $x_{(k+2)}$ does not change which endpoints are in I^- and which are in I^+ – and thus, does not change the corresponding value of the variance.

So:

- for $k < k_{\min}$, we have $r_k > x_{(k+1)}$,
- for $k > k_{\min}$, we have $r_k < x_{(k)}$, and
- for $k = k_{\min}$ (or, to be more precise, for $l_0 \leq k \leq k_0$), we have $x_{(k)} \leq r_k \leq x_{(k+1)}$.

Hence:

- if $r_k < x_{(k)}$, then we cannot have $k < k_{\min}$ and $k = k_{\min}$, hence $k > k_{\min}$;
- if $r_k > x_{(k+1)}$, then we cannot have $k > k_{\min}$ and $k = k_{\min}$, hence $k < k_{\min}$;
- if $x_{(k)} \leq r_k \leq x_{(k+1)}$, then we cannot have $k < k_{\min}$ and $k > k_{\min}$, hence $k = k_{\min}$.

Thus, the above algorithm finds the correct value of k_{\min} and thence, the correct value of \underline{V} .

Proof that the new algorithm for computing \underline{S} requires linear time is similar to the previous proofs.

Proof that under the subset property, the new algorithm always computes \underline{S} . Due to the subset property, we can sort the intervals in lexicographic order, in which case their lower endpoints \underline{p}_i , their upper endpoints \bar{p}_i , and their midpoints \tilde{p}_i are also sorted: $\underline{p}_i \leq \underline{p}_{i+1}$, $\bar{p}_i \leq \bar{p}_{i+1}$, and $\tilde{p}_i \leq \tilde{p}_{i+1}$. Let us thus assume that the intervals are thus sorted.

Let us now show that it is sufficient to consider monotonic optimal tuples p_1, \dots, p_n , for which $p_i \leq p_{i+1}$ for all i . Indeed, if $p_i > p_{i+1}$, then, since $p_i \leq \bar{p}_i \leq \bar{p}_{i+1}$ and $p_i > p_{i+1} \geq \underline{p}_{i+1}$, we have $p_i \in [\underline{p}_{i+1}, \bar{p}_{i+1}]$ and similarly $p_{i+1} \in [\underline{p}_i, \bar{p}_i]$. Thus, we can swap the values p_i and p_{i+1} without changing the value of S . We can repeat this swap as many times as necessary until we get a monotonic tuple that has the exact same value $S = \underline{S}$.

Let us now show that in the optimal tuple, at most one p_i can be inside the corresponding interval. Indeed, if we have two values p_j and p_k strictly inside their intervals, then for an arbitrary small Δ , replacing p_j with $p_j - \Delta$ and p_k with $p_k + \Delta = p_j + \Delta$ should not increase the resulting entropy. This is only possible when the derivative of the resulting expression w.r.t. Δ is 0, i.e., when $p_j = p_k$.

Now, for $p_j - \Delta = p_j - \Delta$ and $p_k + \Delta = p_j + \Delta$, the function S should have a minimum at $\Delta = 0$ and thus, its second derivative relative to Δ should be non-negative. However, an explicit computation shows that this derivative is negative. Thus, our assumption is false, and at most one p_j can be inside the corresponding interval.

Since the values \underline{p}_i are sorted by i , and the values \bar{p}_j are sorted by j , we can now conclude that:

- if $p_j = \underline{p}_j$ and $p_m > \underline{p}_m$, then $p_j \leq p_m$; and
- if $p_m = \bar{p}_m$ and $p_j < \bar{p}_j$, then $p_m \geq p_j$.

Thus, each value $p_j = \underline{p}_j$ precedes all the values $p_m = \bar{p}_m$, and the only value p_i which is strictly inside the corresponding interval lies in between these values. Thus, in a monotonic optimal tuple p_1, \dots, p_n , the first elements are equal to \underline{p}_j , then we may have one element which is strictly inside its interval, and then we have values $p_m = \bar{p}_m$.

For the resulting vector $p = (\underline{p}_1, \dots, \underline{p}_{k-1}, p_k, \bar{p}_{k+1}, \dots, \bar{p}_n)$, with $\underline{p}_k \leq p_k \leq \bar{p}_k$, the condition $\sum_{i=1}^n p_i = 1$ implies that $\Sigma_k \leq 1 \leq \Sigma_{k-1}$, where $\Sigma_k \stackrel{\text{def}}{=} \sum_{i=1}^k \underline{p}_i + \sum_{j=k+1}^n \bar{p}_j$. When we go from Σ_k to Σ_{k+1} , we replace a larger value \bar{p}_{k+1} with a smaller value \underline{p}_{k+1} , hence $\Sigma_k > \Sigma_{k+1}$. Thus, there has to be exactly one k_{\max} for which $\Sigma_k \leq 1 \leq \Sigma_{k-1}$.

So, if we have $\Sigma_m > 1$, this means that the value k_{\max} corresponding to the minimum of S is $> m$; hence for all the indices $\leq m$, we already know that in the optimal vector p , $p_i = \underline{p}_i$. Thus, these indices can be added to the set I^- .

If $\Sigma_{m-1} (= \Sigma_m + 2\Delta_m) < 1$, this means that the value k_{\min} corresponding to the minimum of S is $< m$; hence for all the indices $\geq m$, we already know that in the optimal vector p , $p_j = \bar{p}_j$. Thus, these indices can be added to the set I^+ .

Finally, if $\Sigma_m \leq 1 \leq \Sigma_{m-1}$, then this m is where the minimum of S is attained.

The algorithm has been justified.

Proof that the new algorithm for computing \bar{S} requires linear time is similar to the previous proofs.

Proof that the new algorithm always computes \bar{S} . It is known [8, 9, 16] that if we sort all $2n$ endpoints into a sequence $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(2n)}$, then for some $k = k_{\max}$, the maximum \bar{S} is attained for the vector p for which:

- for all indices j for which $\bar{p}_j \leq p_{(k)}$, we have $p_j = \bar{p}_j$;
- for all indices i for which $\underline{p}_i \geq p_{(k+1)}$, we have $p_i = \underline{p}_i$;

- for all other indices, we have $p_i = \text{const}$; since $\sum_{i=1}^n p_i = 1$, we conclude that this constant is equal to

$$r_k \stackrel{\text{def}}{=} \frac{1 - E_k}{n - N_k}, \text{ where}$$

$$E_k = \sum_{j: \bar{p}_j \leq p^{(k)}} \bar{p}_j + \sum_{i: \underline{p}_i \geq p^{(k+1)}} \underline{p}_i;$$

$$N_k = \#\{j : \bar{p}_j \leq p^{(k)}\} + \#\{i : \underline{p}_i \geq p^{(k+1)}\}.$$

It can also be proven that for the optimal k , we have $r_k \in [p^{(k)}, p^{(k+1)}]$. These facts can be proven, e.g., by the same analysis (adding Δp to one value p_j and subtracting Δp from another value p_k) as in our above analysis of \mathcal{S} .

Let us first prove that if $r_k = \frac{1 - E_k}{n - N_k} \leq p^{(k+1)}$, then the similar inequality $r_{k+1} = \frac{1 - E_{k+1}}{n - N_{k+1}} \leq p^{(k+2)}$

holds for the next value k . Indeed, the given inequality $\frac{1 - E_k}{n - N_k} \leq p^{(k+1)}$ is equivalent to $1 - E_k \leq (n - N_k) \cdot p^{(k+1)}$.

The only difference between the sums $E_k = \sum_{j: \bar{p}_j \leq p^{(k)}} \bar{p}_j + \sum_{i: \underline{p}_i \geq p^{(k+1)}} \underline{p}_i$ and $E_{k+1} = \sum_{j: \bar{p}_j \leq p^{(k+1)}} \bar{p}_j + \sum_{i: \underline{p}_i \geq p^{(k+2)}} \underline{p}_i$ is that:

- some terms equal to $p^{(k+1)}$ may be added (if there are j for which $\bar{p}_j = p^{(k+1)}$), and
- some other terms equal to $p^{(k+1)}$ may be subtracted (if there are i for which $\underline{p}_i = p^{(k+1)}$).

In general, $E_{k+1} = E_k + c_k \cdot p^{(k+1)}$ for some integer c_k (positive, negative, or zero), and $N_{k+1} = N_k + c_k$. Subtracting $c_k \cdot p^{(k+1)}$ from both sides of the given inequality $1 - E_k \leq (n - N_k) \cdot p^{(k+1)}$, we conclude that $1 - E_{k+1} \leq (n - N_{k+1}) \cdot p^{(k+1)}$, i.e., that $r_{k+1} = \frac{1 - E_{k+1}}{n - N_{k+1}} \leq p^{(k+1)}$. Since the sequence $p^{(k)}$ is sorted, we thus conclude that $p^{(k+1)} \leq p^{(k+2)}$ and hence $r_{k+1} \leq p^{(k+2)}$.

So, if the inequality $r_k \leq p^{(k+1)}$ holds for some k , it holds for all larger values of k as well. Thus, this inequality holds for all k after a certain value l_0 .

Similarly, we can prove that if the inequality $r_k \geq p^{(k)}$ holds for some k , then it holds for $k - 1$ as well – since the only difference between E_k and E_{k-1} consists of adding and/or subtracting some values $p^{(k)}$. So, if the inequality $r_k \geq p^{(k)}$ holds for some k , it holds for all smaller values of k as well. Thus, this inequality holds for all k until a certain value k_0 .

Similarly to the proof about \underline{V} , we can prove that if there are several values $k = l_0, l_0 + 1, \dots, k_0$ for which both inequalities hold $p^{(k)} \leq r_k \leq p^{(k+1)}$, then for these k , the entropy has exactly the same value.

So:

- for $k < k_{\max}$, we have $r_k > p^{(k+1)}$,
- for $k > k_{\max}$, we have $r_k < p^{(k)}$, and
- for $k = k_{\max}$ (or, to be more precise, for $l_0 \leq k \leq k_0$), we have $p^{(k)} \leq r_k \leq p^{(k+1)}$.

Hence:

- if $r_k < p^{(k)}$, then we cannot have $k < k_{\max}$ and $k = k_{\max}$, hence $k > k_{\max}$;
- if $r_k > p^{(k+1)}$, then we cannot have $k > k_{\max}$ and $k = k_{\max}$, hence $k < k_{\max}$;
- if $p^{(k)} \leq r_k \leq p^{(k+1)}$, then we cannot have $k < k_{\min}$ and $k > k_{\min}$, hence $k = k_{\max}$.

Thus, the above algorithm finds the correct value of k_{\max} and thence, the correct value of \bar{S} .