

Towards Adding Probabilities and Correlations to Interval Computations

Daniel Berleant
Department of Information Science
University of Arkansas at Little Rock
Little Rock, Arkansas 72204, USA
jberleant@ualr.edu

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

ABSTRACT

The traditional engineering approach to error estimation assumes that we know the probabilities of different values of measurement error $\Delta x_i \stackrel{\text{def}}{=} \tilde{x}_i - x_i$. Yet in many practical situations, we only know the upper bound Δ_i for this error. Hence after the measurement, the only information that we have about x_i is that it belongs to the interval $\mathbf{x}_i \stackrel{\text{def}}{=} [\tilde{x}_i - \Delta_i, \tilde{x}_i + \Delta_i]$. In this case, we have a classic *interval computations* problem: find the narrowest possible interval \mathbf{y} enclosing all possible values of the result $y = f(x_1, \dots, x_n)$ when $x_i \in \mathbf{x}_i$. In this paper, we generalize the preceding case by discussing what to do when, in addition to the bounds Δ_i , we permit partial information about the probabilities of different values of Δx_i and their correlations.

Categories and Subject Descriptors

F.2.1 [Theory of Computation]: Analysis of Algorithms and Problem Complexity—*Numerical Algorithms and Problems*; G.1.0 [Mathematics of Computing]: Numerical Analysis—*Error analysis*; G.4 [Mathematics of Computing]: Mathematical Software—*Algorithm design and analysis*

1. FORMULATION OF THE PROBLEM

Why indirect measurements? In many real-life situations, we are interested in the value of a physical quantity y that is difficult or impossible to measure directly. Examples of such quantities are the distance to a star and the amount of oil in a given well. Since we cannot measure y directly, a natural strategy is to measure y *indirectly*. Specifically, we find some easier-to-measure quantities x_1, \dots, x_n which are related to y by a known relation $y = f(x_1, \dots, x_n)$. To estimate y , we first obtain measurements $\tilde{x}_1, \dots, \tilde{x}_n$ of the quantities x_1, \dots, x_n , and then compute an estimate for y of $\tilde{y} = f(\tilde{x}_1, \dots, \tilde{x}_n)$.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
Copyright 2007 ACM ...\$5.00.

Why interval computations? Measurements are never 100% accurate, so the actual value x_i of measured quantity i can differ from the measurement result \tilde{x}_i . Because of these *measurement errors* $\Delta x_i \stackrel{\text{def}}{=} \tilde{x}_i - x_i$, the result $\tilde{y} = f(\tilde{x}_1, \dots, \tilde{x}_n)$ is, in general, different from the actual value $y = f(x_1, \dots, x_n)$ of the desired quantity y [8].

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

What do we know about the errors Δx_i of direct measurements? First, the manufacturer of the measuring instrument may supply us with an upper bound Δ_i on the measurement error. In this case, once we perform a measurement and get a measurement result \tilde{x}_i , we know that the actual (unknown) value x_i of the measured quantity is in the interval $\mathbf{x}_i = [\underline{x}_i, \bar{x}_i]$, where $\underline{x}_i = \tilde{x}_i - \Delta_i$ and $\bar{x}_i = \tilde{x}_i + \Delta_i$.

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

In this case, after performing a measurement and getting a measurement result \tilde{x}_i , the only information that we have about the actual value x_i of the measured quantity is that it belongs to the interval $\mathbf{x}_i = [\tilde{x}_i - \Delta_i, \tilde{x}_i + \Delta_i]$.¹ In such situations, the only information that we have about the (unknown) actual value of $y = f(x_1, \dots, x_n)$ is that y belongs to the range $\mathbf{y} = [\underline{y}, \bar{y}]$ of the function f over the box $\mathbf{x}_1 \times \dots \times \mathbf{x}_n$:

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

The process of computing this interval range based on the input intervals \mathbf{x}_i is part of *interval computations*; see, e.g., [3].

Interval computations techniques: brief reminder. Historically what is often called the “straightforward” method was the first for estimating the desired range of a function. This method is based on the fact that inside the computer, every algorithm for processing real numbers is implemented as a sequence of elementary operations $a + b$, $a - b$, $a \cdot b$, and a/b ; usually, a/b is computed as $a \cdot (1/b)$, making $a + b$, $a - b$, $a \cdot b$, and $1/a$ sufficient. For each of these elementary operations $f(a, b)$, if we know the intervals \mathbf{a} and \mathbf{b} for a and b , we can compute the exact range $f(\mathbf{a}, \mathbf{b})$. The corresponding formulas form the so-called *in-*

¹We use the convention of bold, non-italic symbols for naming intervals.

terval arithmetic:

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

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

$$[\underline{a}, \bar{a}] \cdot [\underline{b}, \bar{b}] =$$

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

$$1/[\underline{a}, \bar{a}] = [1/\bar{a}, 1/\underline{a}] \text{ if } 0 \notin [\underline{a}, \bar{a}].$$

In straightforward interval computations, we replace each floating point operation in the program f by the corresponding interval operation. It is known that, as a result, we get an enclosure $\mathbf{Y} \supseteq \mathbf{y}$ of the desired range.

In some cases, $\mathbf{Y} = \mathbf{y}$. In more complex cases, the enclosure has excess width ($\mathbf{Y} \supset \mathbf{y}$). There exist more sophisticated techniques for producing narrower enclosures, e.g., centered form methods [3]. However, for each of these techniques, there are cases when we still get excess width. Reason: it is known (see, e.g., [6]), that the problem of computing the exact range is NP-hard even for polynomial functions $f(x_1, \dots, x_n)$ (indeed, even for quadratic functions f).

Motivating practical problem. In some practical situations, in addition to lower and upper bounds on each random variable x_i , we know bounds $\mathbf{E}_i = [\underline{E}_i, \bar{E}_i]$ on its mean E_i .

Indeed, in measurement practice (e.g. [8]), the overall measurement error Δx is usually represented as a sum of two components: a *systematic* error component $\Delta_s x$ which is defined as the expected value $E[\Delta x]$, and a *random* error component $\Delta_r x$ which is defined as the difference between overall measurement error Δx and the systematic error component $\Delta_s x$: $\Delta_r x \stackrel{\text{def}}{=} \Delta x - \Delta_s x$. In addition to an upper bound Δ on the magnitude of overall measurement errors, the manufacturers of a measuring instrument often provide an upper bound Δ_s on the magnitude of the systematic error component: $|\Delta_s x| \leq \Delta_s$.

When this additional information is given, then, after obtaining a measurement result \tilde{x} , we not only have the information that the actual value x of the measured quantity belongs to the interval $\mathbf{x} = [\tilde{x} - \Delta, \tilde{x} + \Delta]$, but we can also conclude that the expected value $E[x]$ of $x = \tilde{x} - \Delta x$ (which is $E[x] = \tilde{x} - E[\Delta x] = \tilde{x} - \Delta_s x$) belongs to the interval $[\tilde{x} - \Delta_s, \tilde{x} + \Delta_s]$.

If we have this information for every x_i , then, in addition to the interval \mathbf{y} of possible values of y , we can also know the interval of possible values of $E[y]$. This additional interval will, we hypothesized, provide us with information on how repeated measurements can improve the accuracy of this indirect measurement. Thus, we arrive at the following problem.

New problem in precise terms. Given an algorithm computing a function $f(x_1, \dots, x_n)$ from \mathbb{R}^n to \mathbb{R} , and values $\underline{x}_1, \bar{x}_1, \dots, \underline{x}_n, \bar{x}_n, \underline{E}_1, \bar{E}_1, \dots, \underline{E}_n, \bar{E}_n$, we want to find

$$\underline{E} \stackrel{\text{def}}{=} \min\{E[f(x_1, \dots, x_n)] : \text{all distributions of}$$

$$(x_1, \dots, x_n) \text{ for which } x_1 \in [\underline{x}_1, \bar{x}_1], \dots, x_n \in [\underline{x}_n, \bar{x}_n],$$

$$E[x_1] \in [\underline{E}_1, \bar{E}_1], \dots, E[x_n] \in [\underline{E}_n, \bar{E}_n]\};$$

and \bar{E} which is the maximum of $E[f(x_1, \dots, x_n)]$ for all such distributions.

In addition to considering all possible distributions, we can also consider the case when all the variables x_i are independent, or, more generally, when we know the correlations among the x_i .

2. WHAT IS KNOWN

Extending interval arithmetic to handle expectations. The main idea behind standard interval computations can be applied here as well. First we find out how to solve the problem when $n = 2$ and $f(x_1, x_2)$ is one of the standard arithmetic operations. Then, once we have an arbitrary algorithm $f(x_1, \dots, x_n)$, we parse it and replace each elementary operation on real numbers with the corresponding operation on quadruples $(\underline{x}, \underline{E}, \bar{E}, \bar{x})$.

To implement this idea, we must therefore know how to solve the above problem for elementary operations.

For *addition*, the answer is straightforward: $E[x_1 + x_2] = E[x_1] + E[x_2]$. So, if $y = x_1 + x_2$, the only possible value for $E = E[y]$ is $E = E_1 + E_2$. This value does not depend on whether we have correlation or whether we have any information about the correlation. Thus, $\mathbf{E} = \mathbf{E}_1 + \mathbf{E}_2$.

Similarly, the answer is straightforward for *subtraction*: if $y = x_1 - x_2$, there is only one possible value for $E = E[y]$: the value $E = E_1 - E_2$. Thus, $\mathbf{E} = \mathbf{E}_1 - \mathbf{E}_2$.

For *multiplication*, if the variables x_1 and x_2 are independent, then $E[x_1 \cdot x_2] = E[x_1] \cdot E[x_2]$. Hence, if $y = x_1 \cdot x_2$ and x_1 and x_2 are independent, there is only one possible value for $E = E[y]$: the value $E = E_1 \cdot E_2$; hence $\mathbf{E} = \mathbf{E}_1 \cdot \mathbf{E}_2$.

The only non-trivial case is the case of multiplication in the presence of possible correlation. When we know the exact values of E_1 and E_2 , the solution to the above problem is known [4]:

THEOREM 1. *If $y = x_1 \cdot x_2$, and we have no information about the correlation, then the range $[\underline{E}, \bar{E}]$ of $E[x_1 \cdot x_2]$ is $[E_{\min}, E_{\max}]$, where $p_i \stackrel{\text{def}}{=} (E_i - \underline{x}_i)/(\bar{x}_i - \underline{x}_i)$, and:*

$$E_{\min} \stackrel{\text{def}}{=} \max(p_1 + p_2 - 1, 0) \cdot \bar{x}_1 \cdot \bar{x}_2 +$$

$$\min(p_1, 1 - p_2) \cdot \bar{x}_1 \cdot \underline{x}_2 + \min(1 - p_1, p_2) \cdot \underline{x}_1 \cdot \bar{x}_2 + \quad (1)$$

$$\max(1 - p_1 - p_2, 0) \cdot \underline{x}_1 \cdot \underline{x}_2;$$

$$E_{\max} \stackrel{\text{def}}{=} \min(p_1, p_2) \cdot \bar{x}_1 \cdot \bar{x}_2 +$$

$$\max(p_1 - p_2, 0) \cdot \bar{x}_1 \cdot \underline{x}_2 + \max(p_2 - p_1, 0) \cdot \underline{x}_1 \cdot \bar{x}_2 + \quad (2)$$

$$\min(1 - p_1, 1 - p_2) \cdot \underline{x}_1 \cdot \underline{x}_2.$$

Comment. In this case, $\mathbf{E} = [E_{\min}, E_{\max}]$. In the following text, we will use the expressions (1) and (2) to describe the ranges of E for other cases, when the expression for the range $\mathbf{E} = [\underline{E}, \bar{E}]$ is different from the above expression $[E_{\min}, E_{\max}]$.

For the *inverse* $y = 1/x_1$, a finite range is possible only when $0 \notin \mathbf{x}_1$. Without loss of generality, we can consider the case when $0 < \underline{x}_1$. In this case, we have the following bound [4]:

THEOREM 2. *For the inverse $y = 1/x_1$, the range of possible values of E is $\mathbf{E} = [1/E_1, p_1/\bar{x}_1 + (1 - p_1)/\underline{x}_1]$.*

(Here p_1 denotes the same value as in Theorem 1.)

Taking correlation into account. As we have seen, for elementary arithmetic operations other than multiplication, the range of the result's expectation is uniquely determined by the ranges of the input expectations. For multiplication, the range of $E[x_1 \cdot x_2]$ depends on both the ranges of $E[x_i]$ and the correlation between the x_i .

For multiplication, we know the bounds on $E[x_1 \cdot x_2]$ for two cases: when x_1 and x_2 are independent, and when we have no information about their correlation. In reality, we may have partial information about the correlation. For example, we may know the exact value ρ of the correlation

$$\rho(x_1, x_2) \stackrel{\text{def}}{=} \frac{E[x_1 \cdot x_2] - E_1 \cdot E_2}{\sigma_1 \cdot \sigma_2} \quad (3)$$

(where σ_i is the standard deviation of x_i). Or more generally we might have an interval $[\underline{\rho}, \bar{\rho}]$ of possible values of ρ .

Analytical expressions are desirable. In [1], a linear programming-based numerical method is described for computing the ranges of binary functions under constraints on the correlation of its arguments. For example, this method can be applied to the problem of estimating the range of $E[x_1 \cdot x_2]$ under known correlation.

In the cases of independence and unknown correlation, there are explicit analytical expressions for the range of $E[x_1 \cdot x_2]$. In general, analytical expressions are much faster to compute than numerical methods. In this paper, we provide analytical expressions for the correlation case as well.

3. MAIN RESULTS

Preliminaries. Our objective is, given the intervals $[\underline{x}_1, \bar{x}_1]$, $[\underline{x}_2, \bar{x}_2]$, the values $E_1 = E[x_1]$, $E_2 = E[x_2]$, and $\rho = \rho(x_1, x_2)$, to find the range $[\underline{E}, \bar{E}]$ of possible values of $E[x_1 \cdot x_2]$.

Before we derive an expression for the general situation, let us identify the quantitative values for Pearson correlation coefficient ρ corresponding to the known cases – independence and unknown correlation. For the former case, $\rho = 0$. For the latter, according to [4] both E_{\min} and E_{\max} are attained when each of the variables x_i has a 2-point (2-impulse) marginal distribution: $p(x_i = \bar{x}_i) = p_i$ and $p(x_i = \underline{x}_i) = 1 - p_i$. (Probability p_i is uniquely determined by expected value $E[x_i]$.) For this marginal distribution,

$$\sigma^2[x_i] = E[(x_i - E_i)^2] = p_i \cdot (\bar{x}_i - E_i)^2 + (1 - p_i) \cdot (E_i - \underline{x}_i)^2.$$

Since $p_i = (E_i - \underline{x}_i) / (\bar{x}_i - \underline{x}_i)$, algebraic manipulation yields

$$\sigma^2[x_i] = (\bar{x}_i - E_i) \cdot (E_i - \underline{x}_i).$$

Thus, using eq. (3), the correlation coefficients ρ_{\min} and ρ_{\max} corresponding to these extreme distributions are equal to $\rho_{\min} = \frac{E_{\min} - E_1 \cdot E_2}{\sigma}$ and $\rho_{\max} = \frac{E_{\max} - E_1 \cdot E_2}{\sigma}$, where

$$\sigma \stackrel{\text{def}}{=} \sigma_1 \cdot \sigma_2 = \sigma[x_1] \cdot \sigma[x_2] =$$

$$\sqrt{(\bar{x}_1 - E_1) \cdot (E_1 - \underline{x}_1)} \cdot \sqrt{(\bar{x}_2 - E_2) \cdot (E_2 - \underline{x}_2)}.$$

Case of exactly known non-zero correlation. The negative value ρ_{\min} corresponds to the smallest possible value E_{\min} of $E[x_1 \cdot x_2]$, and the positive value ρ_{\max} corresponds to the largest possible value E_{\max} . Because the corresponding

analyses are limited to the extremes, it is therefore desirable to extend results to include intermediate values of ρ .

THEOREM 3. Let $[\underline{x}_1, \bar{x}_1]$ and $[\underline{x}_2, \bar{x}_2]$ be given intervals, $E_1 \in [\underline{x}_1, \bar{x}_1]$ and $E_2 \in [\underline{x}_2, \bar{x}_2]$ be given numbers, and ρ be a number from the interval $[\rho_{\min}, \rho_{\max}]$. Then the closure $[\underline{E}, \bar{E}]$ of the range of possible values $E[x_1, x_2]$ for all possible distributions for which:

- x_1 is located in $[\underline{x}_1, \bar{x}_1]$, and x_2 is located in $[\underline{x}_2, \bar{x}_2]$;
- $E[x_1] = E_1$, and $E[x_2] = E_2$; and
- $\rho[x_1, x_2] = \rho$,

is

- for $\rho \geq 0$: $[E_1 \cdot E_2, E_1 \cdot E_2 + \rho \cdot \sigma]$;
- for $\rho \leq 0$: $[E_1 \cdot E_2 + \rho \cdot \sigma, E_1 \cdot E_2]$.

Comment. The need for closure comes from the fact that ρ is only defined when $\sigma_i > 0$. Thus, e.g., for $\rho > 0$, eq. (3) implies $E[x_1 \cdot x_2] > E[x_1] \cdot E[x_2]$. So, under the standard definition of (Pearson) correlation, the lower endpoint $E_1 \cdot E_2$ might be unattainable.

If we instead define a distribution with correlation ρ as a distribution for which

$$E[x_1 \cdot x_2] = E[x_1] \cdot E[x_2] + \rho \cdot \sigma[x_1] \cdot \sigma[x_2],$$

then the degenerate distribution $x_1 \equiv E_1$, $x_2 \equiv E_2$, with $\sigma[x_1] = \sigma[x_2] = 0$, is a distribution with a given ρ for which $E[x_1 \cdot x_2] = E_1 \cdot E_2$. Under this alternative definition, closure is not needed.

Proof. When $\rho = 0$, then, by definition of the correlation, $E[x_1 \cdot x_2] = E_1 \cdot E_2$. So, it is sufficient to consider values of $\rho \neq 0$. In this proof, we will only consider the case $\rho > 0$; the case $\rho < 0$ is similar.

We first prove that the value $E[x_1 \cdot x_2]$ always belongs to the interval $[E_1 \cdot E_2, E_1 \cdot E_2 + \rho \cdot \sigma]$. $E_1 \cdot E_2$ is the lower bound because, since $\rho > 0$, we have $E[x_1 \cdot x_2] = E_1 \cdot E_2 + \rho \cdot \sigma[x_1] \cdot \sigma[x_2] > E_1 \cdot E_2$.

To prove the upper bound, we show that for each x_i , $\sigma^2[x_i] \leq (E_i - \underline{x}_i) \cdot (\bar{x}_i - E_i)$. Let us first consider discrete distributions that take values $x_i^{(j)} \in [\underline{x}_i, \bar{x}_i]$ ($1 \leq j \leq N$) with probabilities $p^{(j)} \geq 0$ such that $\sum_{j=1}^N p^{(j)} = 1$. For such distributions, the constraint $E[x_i] = E_i$ takes the form $\sum_{j=1}^N p^{(j)} \cdot x_i^{(j)} = E_i$. Under these constraints, let us find the largest possible value of

$$\sigma^2[x_i] = E[x_i^2] - E_i^2 = \sum_{j=1}^N p^{(j)} \cdot (x_i^{(j)})^2 - E_i^2.$$

In terms of the unknown probabilities $p_i^{(j)}$, we are minimizing a linear function under linear constraints (equalities and inequalities). Geometrically, the set of all points that satisfy several linear constraints is a polytope. It is well known that to find the minimum of a linear function on a polytope, it is sufficient to consider its vertices (this is the idea behind linear programming). In algebraic terms, a vertex can be characterized by the fact that for N variables, N of the original constraints are equalities. Thus, in our case,

all but two probabilities $p_i^{(j)}$ must be equal to 0, i.e., the distribution must be located at two points x_i^- and x_i^+ . Since the mean is E_i , we these values must be on different sides of E_i . Without losing generality, we can thus assume that $x_i^- \leq E_i \leq x_i^+$.

We have already mentioned that for 2-point distributions, once the points x_i^- and x_i^+ are fixed, the condition that the mean equals E_i uniquely determines the probabilities, and the resulting variance is $(x_i^+ - E_i) \cdot (E_i - x_i^-)$. When $x_i^+ \leq \bar{x}_i$ and $x_i^- \geq \underline{x}_i$, the largest value of this product is attained when x_i^+ attains its largest possible value \bar{x}_i , and x_i^- attains its smallest possible value \underline{x}_i . Thus, for discrete distributions, $\sigma^2[x_i] \leq (\bar{x}_i - E_i) \cdot (E_i - \underline{x}_i)$.

An arbitrary distribution can be approximated by discrete ones to arbitrary accuracy (in weak topology), so this inequality is true for all distributions. Thus, $\sigma[x_1] \cdot \sigma[x_2] \leq \sigma$, and the equality $E[x_1 \cdot x_2] = E_1 \cdot E_2 + \rho \cdot \sigma[x_2] \cdot \sigma[x_2]$ implies that $E[x_1 \cdot x_2] \leq E_1 \cdot E_2 + \rho \cdot \sigma$.

We now prove that both endpoints are exact. For every $\varepsilon > 0$, if we take a distribution in which each x_i is located in the ε -vicinity of E_i , then $x_1 \cdot x_2$ (and hence $E[x_1 \cdot x_2]$) is located in the close vicinity of $E_1 \cdot E_2$. When $\varepsilon \rightarrow 0$, we conclude that $E[x_1 \cdot x_2]$ can be arbitrarily close to $E_1 \cdot E_2$, so the lower endpoint is indeed exact.

To complete the proof, we next show that the upper endpoint $E_1 \cdot E_2 + \rho \cdot \sigma$ is attainable, and thus also exact. Indeed, as we have mentioned, the largest possible value E_{\max} is attained for a joint distribution in which both marginal distributions are 2-point ones, located on the endpoints of the corresponding interval $[\underline{x}_i, \bar{x}_i]$, and that for such distributions, $\sigma^2[x_i] = (\bar{x}_i - E_i) \cdot (E_i - \underline{x}_i)$. In general, distributions with such marginals are located at 4 vertices of the rectangle $[\underline{x}_1, \bar{x}_1] \times [\underline{x}_2, \bar{x}_2]$. The set of such distributions is determined by linear constraints and is, thus, connected. Along this set, the correlation ranges from 0 to the value ρ_{\max} . Since $\rho \in [0, \rho_{\max}]$ and correlation continuously depends on the probabilities, there exists an intermediate value of these probabilities where the correlation exactly equals the given value ρ .

The theorem is proven.

Case of correlation known with interval uncertainty.

We can handle the case of an interval $[\underline{\rho}, \bar{\rho}]$ of possible values for ρ instead of an exact value of ρ by simply combining the intervals from Theorem 3 and using the fact that the corresponding formulas monotonically depend on ρ .

THEOREM 4. *Let $[\underline{x}_1, \bar{x}_1]$ and $[\underline{x}_2, \bar{x}_2]$ be given intervals, $E_1 \in [\underline{x}_1, \bar{x}_1]$ and $E_2 \in [\underline{x}_1, \bar{x}_1]$ be given numbers, and $[\underline{\rho}, \bar{\rho}]$ be a subinterval of the interval $[\rho_{\min}, \rho_{\max}]$. Then the closure $[\underline{E}, \bar{E}]$ of the range of possible values $E[x_1, x_2]$ for all possible distributions for which:*

- x_1 is located in $[\underline{x}_1, \bar{x}_1]$, and x_2 is located in $[\underline{x}_2, \bar{x}_2]$;
- $E[x_1] = E_1$, and $E[x_2] = E_2$; and
- $\rho[x_1, x_2] \in [\underline{\rho}, \bar{\rho}]$

equals

- for $0 \leq \underline{\rho}$: $[E_1 \cdot E_2, E_1 \cdot E_2 + \bar{\rho} \cdot \sigma]$;
- for $\bar{\rho} \leq 0$: $[E_1 \cdot E_2 + \underline{\rho} \cdot \sigma, E_1 \cdot E_2]$;
- for $\underline{\rho} \leq 0 \leq \bar{\rho}$: $[E_1 \cdot E_2 + \underline{\rho} \cdot \sigma, E_1 \cdot E_2 + \bar{\rho} \cdot \sigma]$.

4. AUXILIARY RESULTS

Computationally efficient expressions for E_{\min} and E_{\max} .

PROPOSITION 1.

$$E_{\max} = E_1 \cdot E_2 + \min((E_1 - \underline{x}_1) \cdot (\bar{x}_2 - E_2), (\bar{x}_1 - E_1) \cdot (E_2 - \underline{x}_2));$$

$$E_{\min} = E_1 \cdot E_2 - \min((E_1 - \underline{x}_1) \cdot (E_2 - \underline{x}_2), (\bar{x}_1 - E_1) \cdot (\bar{x}_2 - E_2)).$$

Proof. Let us first simplify the expression for E_{\max} from Theorem 1. When $p_1 \leq p_2$, we get

$$E_{\max} = p_1 \cdot \bar{x}_1 \cdot \bar{x}_2 + (p_2 - p_1) \cdot \underline{x}_1 \cdot \bar{x}_2 + (1 - p_2) \cdot \underline{x}_1 \cdot \underline{x}_2 =$$

$$p_1 \cdot (\bar{x}_1 - \underline{x}_1) \cdot \bar{x}_2 + p_2 \cdot \underline{x}_1 \cdot (\bar{x}_2 - \underline{x}_2) + \underline{x}_1 \cdot \underline{x}_2.$$

Substituting the definitions of p_i , we conclude that

$$E_{\max} = (E_1 - \underline{x}_1) \cdot \bar{x}_2 + (E_2 - \underline{x}_2) \cdot \underline{x}_1 + \underline{x}_1 \cdot \underline{x}_2.$$

Opening parentheses, we get

$$E_{\max} = E^{(1)} \stackrel{\text{def}}{=} E_1 \cdot \bar{x}_2 - \underline{x}_1 \cdot \bar{x}_2 + E_2 \cdot \underline{x}_1.$$

By using the symmetry between x_1 and x_2 , we can now conclude that when $p_1 \geq p_2$,

$$E_{\max} = E^{(2)} \stackrel{\text{def}}{=} E_2 \cdot \bar{x}_1 - \bar{x}_1 \cdot \underline{x}_2 + E_1 \cdot \underline{x}_2.$$

The condition $p_1 \leq p_2$ is equivalent to

$$(E_1 - \underline{x}_1) \cdot (\bar{x}_2 - \underline{x}_2) \leq (E_2 - \underline{x}_2) \cdot (\bar{x}_1 - \underline{x}_1),$$

i.e.,

$$E_1 \cdot \bar{x}_2 - E_1 \cdot \underline{x}_2 - \underline{x}_1 \cdot \bar{x}_2 + \underline{x}_1 \cdot \underline{x}_2 \leq E_2 \cdot \bar{x}_1 - E_2 \cdot \underline{x}_1 - \bar{x}_1 \cdot \underline{x}_2 + \underline{x}_1 \cdot \underline{x}_2.$$

Subtracting the common term $\underline{x}_1 \cdot \underline{x}_2$ from both sides and moving terms to other sides, we get an equivalent form of this inequality:

$$E_1 \cdot \bar{x}_2 - \underline{x}_1 \cdot \bar{x}_2 + E_2 \cdot \underline{x}_1 \leq E_2 \cdot \bar{x}_1 - \bar{x}_1 \cdot \underline{x}_2 + E_1 \cdot \underline{x}_2,$$

i.e., $E^{(1)} \leq E^{(2)}$. So, if $p_1 \leq p_2$, i.e., if $E^{(1)} \leq E^{(2)}$, we get $E_{\max} = E^{(1)}$; otherwise, we get $E_{\max} = E^{(2)}$. These two cases can be combined into a single formula $E_{\max} = \min(E^{(1)}, E^{(2)})$, i.e.,

$$E_{\max} = \min(E_1 \cdot \bar{x}_2 - \underline{x}_1 \cdot \bar{x}_2 + E_2 \cdot \underline{x}_1, E_2 \cdot \bar{x}_1 - \bar{x}_1 \cdot \underline{x}_2 + E_1 \cdot \underline{x}_2).$$

By adding $-E_1 \cdot E_2$ to both expressions $E^{(1)}$ and $E^{(2)}$, we get the desired expression for E_{\max} .

Since $E[x_1 \cdot x_2] = -E[(-x_1) \cdot x_2]$, where $-x_1 \in [-\bar{x}_1, \underline{x}_1]$ with $E[-x_1] = -E_1$, we have

$$E_{\min} \stackrel{\text{def}}{=} \min E[x_1 \cdot x_2] = -\max E[(-x_1) \cdot x_2].$$

Hence, the new expression for E_{\max} leads to the desired expression for E_{\min} . The proposition is proven.

Can we propagate correlations through computations?

In straightforward interval computations, we propagate intervals through computations; can we similarly propagate correlations? The following result shows that it is not easy even for addition:

PROPOSITION 2. *If we know that $\rho[x_1, x_2] = \rho$, then the only possible conclusion about $\rho' = \rho[x_1, x_1 + x_2]$ is that $\rho' \in [\rho, 1]$.*

Proof. If we take $x_1 \ll x_2$, we get $\rho' \approx \rho$, and if we take $x_2 \ll x_1$, we get $\rho' \approx 1$. The smaller the corresponding ratio x_1/x_2 or x_2/x_1 , the closer we are, correspondingly, to ρ and to 1.

Let us prove that ρ' cannot be smaller than ρ . Since correlation can be defined in terms of the differences $x_i - E[x_i]$, we can shift both variables to $E[x_i] = 0$ without changing the correlations $\rho[x_1, x_2]$ and $\rho[x_1, x_1 + x_2]$; thus, is it sufficient to prove the desired inequality $\rho' \geq \rho$ for the case when $E[x_i] = 0$. In this case, if we denote $\sigma_i \stackrel{\text{def}}{=} \sigma[x_i]$, we get

$$\rho' = \frac{E[x_1 \cdot (x_1 + x_2)]}{\sigma_1 \cdot \sigma[x_1 + x_2]} = \frac{\sigma_1^2 + E[x_1 \cdot x_2]}{\sigma_1 \cdot \sigma[x_1 + x_2]}.$$

Here, since $E_i = 0$, we have $E[x_1 \cdot x_2] = \rho \cdot \sigma_1 \cdot \sigma_2$. Similarly, $\sigma^2[x_1 + x_2] = E[(x_1 + x_2)^2] = E[x_1^2] + E[x_2^2] + 2 \cdot E[x_1 \cdot x_2] =$

$$\sigma_1^2 + \sigma_2^2 + 2\rho \cdot \sigma_1 \cdot \sigma_2,$$

so the above expression for ρ' takes the form: $\rho' = \frac{\sigma_1 + \rho \cdot \sigma_1 \cdot \sigma_2}{\sigma_1 \cdot \sqrt{\sigma_1^2 + \sigma_2^2 + 2\rho \cdot \sigma_1 \cdot \sigma_2}}$, and the desired inequality $\rho' \geq$

ρ takes the form $\frac{\sigma_1 + \rho \cdot \sigma_2}{\sqrt{\sigma_1^2 + \sigma_2^2 + 2\rho \cdot \sigma_1 \cdot \sigma_2}} \geq \rho$. Multiplying both sides by the denominator, we get the equivalent inequality

$$\sigma_1 + \rho \cdot \sigma_2 \geq \rho \cdot \sqrt{\sigma_1^2 + \sigma_2^2 + 2\rho \cdot \sigma_1 \cdot \sigma_2}. \quad (4)$$

If $\rho \geq 0$, then we can square both sides and get an equivalent inequality

$$\sigma_1^2 + 2\rho \cdot \sigma_1 \cdot \sigma_2 + \rho^2 \cdot \sigma_2^2 \geq \rho^2 \cdot (\sigma_1^2 + \sigma_2^2 + 2\rho \cdot \sigma_1 \cdot \sigma_2).$$

Subtracting $\rho^2 \cdot \sigma_2^2$ from both sides, and moving all the terms to the right-hand side, we get an equivalent inequality

$$\sigma_1^2 \cdot (1 - \rho^2) + 2\rho \cdot \sigma_1 \cdot \sigma_2 \cdot (1 - \rho^2) \geq 0,$$

which is always true for $\rho \geq 0$ (since $\rho \leq 1$).

If $\rho < 0$, the right-hand side of (4) is negative, so we consider two possible cases. The first case is when

$$\sigma_1 + \rho \cdot \sigma_2 \geq 0.$$

Then inequality (4) is automatically true.

The second case is when $\sigma_1 + \rho \cdot \sigma_2 < 0$. In this case, (4) is equivalent to

$$0 < -\sigma_1 + |\rho| \cdot \sigma_2 \leq |\rho| \cdot \sqrt{\sigma_1^2 + \sigma_2^2 - 2|\rho| \cdot \sigma_1 \cdot \sigma_2}.$$

By squaring both sides, we get an equivalent inequality

$$\sigma_1^2 - 2|\rho| \cdot \sigma_1 \cdot \sigma_2 + \rho^2 \cdot \sigma_2^2 \leq \rho^2 \cdot (\sigma_1^2 + \sigma_2^2 - 2|\rho| \cdot \sigma_1 \cdot \sigma_2).$$

Subtracting $\rho^2 \cdot \sigma_2^2$ from both sides, and moving all the terms to the right-hand side, we get an equivalent inequality

$$\sigma_1^2 \cdot (1 - \rho^2) - 2|\rho| \cdot \sigma_1 \cdot \sigma_2 \cdot (1 - \rho^2) \leq 0.$$

Dividing both sides by $\sigma_1 \cdot (1 - \rho^2) > 0$, we get an equivalent inequality $\sigma_1 - 2|\rho| \cdot \sigma_2 \leq 0$. We consider the case when $\sigma_1 - |\rho| \cdot \sigma_2 < 0$, hence $\sigma_1 - 2|\rho| \cdot \sigma_2 \leq \sigma_1 - |\rho| \cdot \sigma_2 < 0$. The inequality is proven.

Since $x_1 - x_2 = x_1 + (-x_2)$, and $\rho[x_1, -x_2] = -\rho[x_1, x_2]$, we have the following corollary:

PROPOSITION 3. *If we know that $\rho[x_1, x_2] = \rho$, then:*

- *the best possible conclusion about $\rho' = \rho[x_1, x_1 - x_2]$ is that $\rho' \in [-\rho, 1]$;*
- *the best possible conclusion about $\rho'' = \rho[x_2, x_1 - x_2]$ is that $\rho'' \in [-1, \rho]$.*

For a *unary* linear function $f(x_1) = a \cdot x_1 + b$, we get $\rho[x_1, f(x_1)] = 1$ for $a > 0$ and $\rho[x_1, f(x_1)] = -1$ for $a < 0$. For non-linear unary functions $f(x_1)$, we can get different intermediate values. As an example, we take $f(x_1) = x_1^2$. Then, $\rho \approx 1$, e.g., for a 2-point distribution located at $a - \varepsilon$ and $a + \varepsilon$ (where $a > 0$ and $\varepsilon \rightarrow 0$) with probability 1/2. $\rho \approx -1$, e.g., for a similar distribution with $a < 0$. We get all possible values from -1 to 1 for intermediate distributions.

5. OPEN PROBLEMS

What if we have a multiple product? For the case of unknown correlation, analytical formulas were obtained in [5].

What if we use different correlation characteristics [9], e.g., the Spearman and Kendall correlations, or copulas [2, 7]?

What about the ranges for $E[\min(x_1, x_2)]$ and $E[\max(x_1, x_2)]$ under a given correlation (for the case of unknown correlation, such ranges were described in [4]).

6. ACKNOWLEDGMENTS

This work was supported in part by NSF grants EAR-0225670 and DMS-0532645 and by Texas Department of Transportation grant No. 0-5453.

7. REFERENCES

- [1] D. Berleant and J. Zhang, "Using Pearson correlation to improve envelopes around the distributions of functions," *Reliable Computing*, 2004, Vol. 10, No. 2, pp. 139–161.
- [2] S. Ferson, *RAMAS RiskCalc: Risk Assessment with Uncertain Numbers*, CRC Press, Boca Raton, Florida, 2002.
- [3] L. Jaulin, M. Keiffer, O. Didrit, E. Walter, *Applied Interval Analysis*, Springer-Verlag, London, 2001.
- [4] V. Kreinovich, "Probabilities, Intervals, What Next? Optimization Problems Related to Extension of Interval Computations to Situations with Partial Information about Probabilities", *Journal of Global Optimization*, 2004, Vol. 29, No. 3, pp. 265–280.
- [5] V. Kreinovich, S. Ferson, and L. Ginzburg, "Exact Upper Bound on the Mean of the Product of Many Random Variables With Known Expectations", *Reliable Computing*, 2003, Vol. 9, No. 6, pp. 441–463.
- [6] V. Kreinovich, A. Lakeyev, J. Rohn, P. Kahl, *Computational complexity and feasibility of data processing and interval computations*, Kluwer, Dordrecht, 1997.
- [7] R. B. Nelsen, *Introduction to Copulas*, Springer Verlag, New York, 1999.
- [8] S. Rabinovich, *Measurement Errors and Uncertainties: Theory and Practice*, Springer-Verlag, New York, 2005.
- [9] D. J. Sheskin, *Handbook of Parametric and Nonparametric Statistical Procedures*, Chapman & Hall/CRC, Boca Raton, Florida, 2004.