

10-2006

How to Efficiently Process Uncertainty Within a Cyberinfrastructure without Sacrificing Privacy and Confidentiality

Luc Longpre

The University of Texas at El Paso, longpre@utep.edu

Vladik Kreinovich

The University of Texas at El Paso, vladik@utep.edu

Follow this and additional works at: https://scholarworks.utep.edu/cs_techrep



Part of the [Computer Engineering Commons](#)

Comments:

UTEP-CS-06-17a.

Published in: Nadia Nedjah, Ajith Abraham, and Luiza de Macedo Mourelle (Eds.), *Computational Intelligence in Information Assurance and Security*, Springer-Verlag, 2007, pp. 155-173.

Recommended Citation

Longpre, Luc and Kreinovich, Vladik, "How to Efficiently Process Uncertainty Within a Cyberinfrastructure without Sacrificing Privacy and Confidentiality" (2006). *Departmental Technical Reports (CS)*. 169.
https://scholarworks.utep.edu/cs_techrep/169

This Article is brought to you for free and open access by the Computer Science at ScholarWorks@UTEP. It has been accepted for inclusion in Departmental Technical Reports (CS) by an authorized administrator of ScholarWorks@UTEP. For more information, please contact lweber@utep.edu.

How to Efficiently Process Uncertainty within a Cyberinfrastructure without Sacrificing Privacy and Confidentiality

Luc Longpré and Vladik Kreinovich

Department of Computer Science, University of Texas at El Paso, 500 W.
University, El Paso, TX 79968, USA, longpre@utep.edu, vladik@utep.edu

In this Chapter, we propose a simple solution to the problem of estimating uncertainty of the results of applying a black-box algorithm – without sacrificing privacy and confidentiality of the algorithm.

1 Cyberinfrastructure and Web Services

1.1 Practical Problem: Need to Combine Geographically Separate Computational Resources

In different knowledge domains in science and engineering, there is a large amount of data stored in different locations, and there are many software tools for processing this data, also implemented at different locations. Users may be interested in different information about this domain.

Sometimes, the information required by the user is already stored in *one of the databases*. For example, if we want to know the geological structure of a certain region in Texas, we can get this information from the geological map stored in Austin. In this case, all we need to do to get an appropriate response to the query is to get this data from the corresponding database.

In other cases, different pieces of the information requested by the user are *stored at different locations*. For example, if we are interested in the geological structure of the Rio Grande Region, then we need to combine data from the geological maps of Texas, New Mexico, and the Mexican state of Chihuahua. In such situations, a correct response to the user's query requires that we access these pieces of information from different databases located at different geographic notations.

In many other situations, the appropriate answer to the user's request requires that we not only collect the relevant data x_1, \dots, x_n , but that we also use some *data processing* algorithms $f(x_1, \dots, x_n)$ to process this data. For example, if we are interested in the large-scale geological structure of a

geographical region, we may also use the gravity measurements from the gravity databases. For that, we need special algorithms to transform the values of gravity at different locations into a map that describes how the density changes with location. The corresponding data processing programs often require a lot of computational resources; as a result, many such programs reside on computers located at supercomputer centers, i.e., on computers which are physically separated from the places where the data is stored.

The need to combine computational resources (data and programs) located at different geographic locations seriously complicates research.

1.2 Centralization of Computational Resources – Traditional Approach to Combining Computational Resources; Its Advantages and Limitations

Traditionally, a widely used way to make these computational resources more accessible was to move all these resources to a *central location*. For example, in the geosciences, the US Geological Survey (USGS) was trying to become a central depository of all relevant geophysical data. However, this centralization requires a large amount of efforts: data are presented in different formats, the existing programs use specific formats, etc. To make the central data depository efficient, it is necessary:

- to reformat all the data,
- to rewrite all the data processing programs – so that they become fully compatible with the selected formats and with each other,
- etc.

The amount of work that is needed for this reformatting and rewriting is so large that none of these central depositories really succeeded in becoming an easy-to-use centralized database.

1.3 Cyberinfrastructure – A More Efficient Approach to Combining Computational Resources

Cyberinfrastructure technique is a new approach that provides the users with the efficient way to submit requests without worrying about the geographic locations of different computational resources – and at the same time avoid centralization with its excessive workloads. The main idea behind this approach is that *we keep all (or at least most) the computational resources*

- *at their current locations,*
- *in their current formats.*

To expedite the use of these resources:

- we supplement the local computational resources with the “metadata”, i.e., with the information about the formats, algorithms, etc.,

- we “wrap up” the programs and databases with auxiliary programs that provide data compatibility into *web services*,

and, in general, we provide a cyberinfrastructure that uses the metadata to automatically combine different computational resources.

For example, if a user is interested in using the gravity data to uncover the geological structure of the Rio Grande region, then the system should automatically:

- get the gravity data from the UTEP and USGS gravity databases,
- convert them to a single format (if necessary),
- forward this data to the program located at San Diego Supercomputer Center, and
- move the results back to the user.

This example is exactly what we are designing under the NSF-sponsored Cyberinfrastructure for the Geosciences (GEON) project; see, e.g., [1, 2, 3, 8, 9, 27, 37, 45, 47, 48, 58, 62, 63]. This is similar to what other cyberinfrastructure projects are trying to achieve.

1.4 What Is Cyberinfrastructure: The Official NSF Definition

According to the final report of the National Science Foundation (NSF) Blue Ribbon Advisory Panel on Cyberinfrastructure, “a new age has dawned in scientific and engineering research, pushed by continuing progress in computing, information, and communication technology, and pulled by the expanding complexity, scope, and scale of today’s challenges. The capacity of this technology has crossed thresholds that now make possible a comprehensive ‘cyberinfrastructure’ on which to build new types of scientific and engineering knowledge environments and organizations and to pursue research in new ways and with increased efficacy.

Such environments and organizations, enabled by cyberinfrastructure, are increasingly required to address national and global priorities, such as understanding global climate change, protecting our natural environment, applying genomics-proteomics to human health, maintaining national security, mastering the world of nanotechnology, and predicting and protecting against natural and human disasters, as well as to address some of our most fundamental intellectual questions such as the formation of the universe and the fundamental character of matter.”

1.5 Web Services: What They Do – A Brief Summary

In different knowledge domains, there is a large amount of data stored in different locations; algorithms for processing this data are also implemented at different locations. Web services – and, more generally, cyberinfrastructure – provide the users with an efficient way to submit requests without worrying

about the geographic locations of different computational resources (databases and programs) – and avoid centralization with its excessive workloads [21]. Web services enable the user to receive the desired data x_1, \dots, x_n and the results $y = f(x_1, \dots, x_n)$ of processing this data.

2 Processing Uncertainty Within a Cyberinfrastructure

2.1 Formulation of the problem

The data x_i usually come from measurements or from experts. Measurements are never 100% accurate; as a result, the measured values \tilde{x}_i are, in general, somewhat different from the actual (unknown) values x_i of the corresponding quantities. Experts can also only provide us with approximate values of the desired quantities.

As a result of this measurement or expert uncertainty, the result $\tilde{y} = f(\tilde{x}_1, \dots, \tilde{x}_n)$ of data processing is, in general, different from the actual value $y = f(x_1, \dots, x_n)$ of the desired quantity. It is desirable to gauge this difference. To do that, we must have some information about the errors of direct measurements.

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

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

These probabilities are often described by a normal distribution, so in standard engineering textbook on measurement, it is usually assumed that the distribution of Δx_i is normal, with 0 average and known standard deviation σ_i ; see, e.g. [20, 46].

In general, we can determine the desired probabilities of different values of Δx_i by comparing the results of measuring with this instrument with the results of measuring the same quantity by a standard (much more accurate) measuring instrument. Since the standard measuring instrument is much more

accurate than the one use, the difference between these two measurement results is practically equal to the measurement error; thus, the empirical distribution of this difference is close to the desired probability distribution for measurement error.

Case of interval uncertainty. There are two cases, however, when in practice, we do not determine the probabilities:

- First is the case of cutting-edge measurements, e.g., measurements in fundamental science. When a Hubble telescope detects the light from a distant galaxy, there is no “standard” (much more accurate) telescope floating nearby that we can use to calibrate the Hubble: the Hubble telescope is the best we have.
- The second case is the case of measurements on the shop floor. In this case, in principle, every sensor can be thoroughly calibrated, but sensor calibration is so costly – usually costing ten times more than the sensor itself – that manufacturers rarely do it.

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

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

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

The process of computing this interval range based on the input intervals \mathbf{x}_i is called *interval computations*; see, e.g., [24, 26].

Case of fuzzy uncertainty. Often, knowledge comes in terms of uncertain expert estimates. In the fuzzy case, to describe this uncertainty, for each value of estimation error Δx_i , we describe the degree $\mu_i(\Delta x_i)$ to which this value is possible.

For each degree of certainty α , we can determine the set of values of Δx_i that are possible with at least this degree of certainty – the α -cut $\{x \mid \mu(x) \geq \alpha\}$ of the original fuzzy set. In most cases, this α -cut is an interval.

Vice versa, if we know α -cuts for every α , then, for each object x , we can determine the degree of possibility that x belongs to the original fuzzy set [7, 30, 38, 40, 41]. A fuzzy set can be thus viewed as a nested family of its α -cuts.

So, if instead of a (crisp) interval \mathbf{x}_i of possible values of the measured quantity, we have a fuzzy set $\mu_i(x)$ of possible values, then we can view this information as a family of nested intervals $\mathbf{x}_i(\alpha)$ – α -cuts of the given fuzzy sets.

2.2 Description of uncertainty: general formulas

In this chapter, we will only consider a typical situation in which the direct measurements and/or expert estimates are accurate enough, so that the resulting approximation errors Δx_i are small, and terms which are quadratic (or of higher order) in Δx_i can be safely neglected. In such situations, the dependence of the desired value $y = f(x_1, \dots, x_n) = f(\tilde{x}_1 - \Delta x_1, \dots, \tilde{x}_n - \Delta x_n)$ on Δx_i can be safely assumed to be linear.

Comment. There are practical situations when the accuracy of the direct measurements is not high enough, and hence, quadratic terms cannot be safely neglected (see, e.g., [26] and references therein). In this case, the problem of error estimation for indirect measurements becomes computationally difficult (NP-hard) even when the function $f(x_1, \dots, x_n)$ is quadratic [34, 60]. However, in most real-life situations, the possibility to ignore quadratic terms is a reasonable assumption, because, e.g., for an error of 1% its square is a negligible 0.01%.

When approximation errors are small, we can simplify the expression for $\Delta y = \tilde{y} - y = f(\tilde{x}_1, \dots, \tilde{x}_n) - f(x_1, \dots, x_n)$ if we expand the function f in Taylor series around the point $(\tilde{x}_1, \dots, \tilde{x}_n)$ and restrict ourselves only to linear terms in this expansion. As a result, we get the expression

$$\Delta y = c_1 \cdot \Delta x_1 + \dots + c_n \cdot \Delta x_n, \quad (1)$$

where by c_i , we denoted the value of the partial derivative $\partial f / \partial x_i$ at the point $(\tilde{x}_1, \dots, \tilde{x}_n)$:

$$c_i = \frac{\partial f}{\partial x_i} \Big|_{(\tilde{x}_1, \dots, \tilde{x}_n)}. \quad (2)$$

Case of probabilistic uncertainty. In the statistical setting, the desired measurement error Δy is a linear combination of independent Gaussian variables Δx_i . Therefore, Δy is also normally distributed, with 0 average and the standard deviation

$$\sigma = \sqrt{c_1^2 \cdot \sigma_1^2 + \dots + c_n^2 \cdot \sigma_n^2}. \quad (3)$$

Comment. A similar formula holds if we *do not* assume that Δx_i are normally distributed: it is sufficient to assume that they are independent variables with 0 average and known standard deviations σ_i .

Case of interval uncertainty. In the interval setting, we do not know the probability of different errors Δx_i ; instead, we only know that $|\Delta x_i| \leq \Delta_i$. In this case, the sum (1) attains its largest possible value if each term $c_i \cdot \Delta x_i$ in this sum attains the largest possible value:

- If $c_i \geq 0$, then this term is a monotonically non-decreasing function of Δx_i , so it attains its largest value at the largest possible value $\Delta x_i = \Delta_i$; the corresponding largest value of this term is $c_i \cdot \Delta_i$.

- If $c_i < 0$, then this term is a decreasing function of Δx_i , so it attains its largest value at the smallest possible value $\Delta x_i = -\Delta_i$; the corresponding largest value of this term is $-c_i \cdot \Delta_i = |c_i| \cdot \Delta_i$.

In both cases, the largest possible value of this term is $|c_i| \cdot \Delta_i$, so, the largest possible value of the sum Δy is

$$\Delta = |c_1| \cdot \Delta_1 + \dots + |c_n| \cdot \Delta_n. \quad (4)$$

Similarly, the smallest possible value of Δy is $-\Delta$.

Hence, the interval of possible values of Δy is $[-\Delta, \Delta]$, with Δ defined by the formula (4).

Case of fuzzy uncertainty. We have already mentioned that if instead of a (crisp) interval \mathbf{x}_i of possible values of the measured quantity, we have a fuzzy set $\mu_i(x)$ of possible values, then we can view this information as a family of nested intervals $\mathbf{x}_i(\alpha)$ – α -cuts of the given fuzzy sets.

Our objective is then to compute the fuzzy number corresponding to this the desired value $y = f(x_1, \dots, x_n)$. In this case, for each level α , to compute the α -cut of this fuzzy number, we can apply interval computations to the α -cuts $\mathbf{x}_i(\alpha)$ of the corresponding fuzzy sets. The resulting nested intervals form the desired fuzzy set for y .

So, e.g., if we want to describe 10 different levels of uncertainty, then we must solve 10 interval computation problems – i.e., apply the formula (4) 10 times.

In many practical situations, there is no need to perform 10 computations. For example, it is often reasonable to assume that all the membership functions $\mu_i(x_i)$ have the same shape and only differ by a scaling parameter, i.e., all have the form $\mu_i(\Delta x_i) = \mu_0(\Delta x_i / \Delta_i)$ for some fixed function $\mu_0(x)$ (e.g., triangular or Gaussian). In this case, as it is well known [14, 15, 16, 25], the membership function for Δy has a similar form $\mu_0(\Delta y / \Delta)$, where Δ is determined by the formula (4).

Comment. These formulas correspond to the case when in the *extension principle* that describe how the uncertainty in Δx_i transforms into the uncertainty in Δy , we interpret “and” as min, i.e., if we consider

$$\mu(\Delta y) = \max_{\Delta x_1, \dots, \Delta x_n} \min(\mu_1(\Delta x_1), \dots, \mu_n(\Delta x_n)), \quad (5)$$

where maximum is taken over all the values $\Delta x_1, \dots, \Delta x_n$ for which the expression (1) for Δy leads to the given value of Δy . In general, we can use a different t-norm to combine the values $\mu_i(\Delta x_i)$. For example, we may use the product and describe the resulting membership function as

$$\mu(\Delta y) = \max_{\Delta x_1, \dots, \Delta x_n} \mu_1(\Delta x_1) \cdot \dots \cdot \mu_n(\Delta x_n). \quad (6)$$

In this case, if we assume that all the membership functions $\mu_i(\Delta x_i)$ are Gaussian, i.e., have the form $\mu_i(\Delta x_i) = \mu_0(\Delta x_i / \sigma_i)$, where $\mu_0(z) = \exp(-z^2)$, then the resulting membership function for Δy is also Gaussian $\mu(\Delta y) = \mu_0(\Delta y / \sigma)$, where σ is determined by the formula (3); see, e.g., [14, 15, 16, 25].

2.3 Error Estimation for the Results of Data Processing: A Precise Computational Formulation of the Problem

As a result of the above analysis, we get the following explicit formulation of the problem: given a function $f(x_1, \dots, x_n)$, n numbers $\tilde{x}_1, \dots, \tilde{x}_n$, and n positive numbers $\sigma_1, \dots, \sigma_n$ (or $\Delta_1, \dots, \Delta_n$), compute the corresponding expression (3) or (4).

Let us describe how this problem is solved now.

2.4 How This Problem Is Solved Now

Textbook case: the function f is given by its analytical expression. If the function f is given by its analytical expression, then we can simply explicitly differentiate it, and get an explicit expression for (3) and (4). This is the case which is typically analyzed in textbooks on measurement theory; see, e.g., [20, 46].

A more complex case: automatic differentiation. In many practical cases, we do not have an explicit analytical expression, we only have an *algorithm* for computing the function $f(x_1, \dots, x_n)$, an algorithm which is too complicated to be expressed as an analytical expression.

When this algorithm is presented in one of the standard programming languages such as Fortran or C, we can let the computer perform an explicit differentiation; for that, we can use one of the existing automatic differentiation tools (see, e.g., [6, 22]). These tools analyze the code of the program for computing $f(x_1, \dots, x_n)$ and, as they perform their analysis, they produce the “differentiation code”, i.e., a program that computes the partial derivatives c_i .

Once we know an algorithm that computes f in time T , automatic differentiation (AD) enables us to compute all partial derivatives in time $\leq 3T$, hence we can compute σ or Δ in time $O(T + n)$.

3 Need for Privacy Makes the Problem More Complex

Privacy situation: description. In cyberinfrastructure, the owners of the program f may not want to disclose its code; instead, they may only allow to use f as a black box.

Real world example. For example, to find places where oil can be found, it is important to know the structure of the Earth crust at different locations. One of the main techniques for determining this structure is the use of seismic data: we use the measured travel-times of the natural (or artificial) seismic signals as they go from their source to their on-surface destination, and then we solve the corresponding *inverse problem* to find the Earth structure; see, e.g., [2, 3, 4, 5, 11, 13, 23, 44, 64].

Because of the importance of the seismic inverse problem, oil companies are heavily investing in developing algorithms for solving such problems. This is a very competitive area of research, so the algorithms are kept confidential. Since most major algorithms have similar efficiency and accuracy, a company is willing to allow other users – researchers or competitors from other companies – to actually use their code. However, companies are very reluctant to let the users have actual access to their code. The reason for this reluctance is that different companies have achieved the similar level of efficiency and accuracy by using different ideas and programming improvements. As a result, a company that unilaterally discloses its code would thus let competitors use its ideas (and thus, improve their code) without getting any benefits in return.

Thus, the companies are allowing to use their code, but only as a black box.

What was known before. The problem of preserving the code’s privacy and confidentiality is a particular case of a general problem of privacy preservation in data processing; see, e.g., [10, 12, 17, 18, 19, 28, 29, 35, 39, 42, 43, 49, 50, 51, 52, 53, 54, 55, 56, 57, 61]. At present, privacy-preserving algorithms are mainly concerned with privacy of *data*; privacy of *code* is a problem for which few results are known.

Related difficulty. If we do not know the code of f , then we cannot apply AD to compute all n partial derivatives $c_i = \frac{\partial f}{\partial x_i}$.

A straightforward method of solving this problem: numerical differentiation. The most straightforward algorithm for solving this problem is to compute the derivatives c_i one-by-one, and then use the corresponding formula (3) or (4) to compute the desired σ . To compute the i -th partial derivative, we change the i -th input x_i to $\tilde{x}_i + h_i$ for some h_i , and leave other inputs unchanged, i.e., we take $\delta_i = h_i$ for this i and $\delta_j = 0$ for all $j \neq i$. Then, we estimate c_i as

$$c_i = \frac{f(\tilde{x}_1, \dots, \tilde{x}_{i-1}, \tilde{x}_i + h_i, \tilde{x}_{i+1}, \dots, \tilde{x}_n) - \tilde{y}}{h_i}. \quad (7)$$

This algorithm is called *numerical differentiation*.

We want the change h_i to be small (so that quadratic terms can be neglected); we already know that changes of the order σ_i are small. So, it is natural to take $h_i = \sigma_i$ (or, correspondingly, $h_i = \Delta_i$). In other words, to compute c_i , we use the following values: $\delta_1 = \dots = \delta_{i-1} = 0$, $\delta_i = \sigma_i$ (or $\delta_i = \Delta_i$), $\delta_{i+1} = \dots = \delta_n = 0$.

Problem: sometimes, numerical differentiation takes too long. Very often, the program f requires a reasonable time to compute (e.g., in the geological applications, computing f may involve solving an inverse problem). In this case, applying the function f is the most time-consuming part of this algorithm. So, the total time that it takes us to compute σ or Δ is (approximately) equal to the running time T for the program f multiplied by the number of times N_f that we call the program f .

For numerical differentiation, $N_f = n$ (we call f n times to compute n partial derivatives). Hence, if the program f takes a long time to compute, and n is huge, then the resulting time $T \cdot n$ (which is $\gg T + n$) may be too long. For example, if we are determining some parameters of an oil well from the geophysical measurements, we may get n in the thousands, and T in minutes. In this case, $T \cdot n$ may take several weeks. This may be OK for a single measurement, but too long if we want more on-line results.

4 Solution for Statistical Setting: Monte-Carlo Simulations

Monte-Carlo simulations: main idea. In the statistical setting, we can use straightforward (Monte-Carlo) simulation, and drastically save the computation time. In this approach, we use a computer-based random number generator to simulate the normally distributed error. A standard normal random number generator usually produces a normal distribution with 0 average and standard deviation 1. So, to simulate a distribution with a standard deviation σ_i , we multiply the result α_i of the standard random number generator by σ_i . In other words, we take $\delta_i = \sigma_i \cdot \alpha_i$.

As a result of N Monte-Carlo simulations, we get N values $c^{(1)} = \mathbf{c} \cdot \boldsymbol{\delta}^{(1)}, \dots, c^{(N)} = \mathbf{c} \cdot \boldsymbol{\delta}^{(N)}$ which are normally distributed with the desired standard deviation σ . So, we can determine σ by using the standard statistical estimate

$$\sigma = \sqrt{\frac{1}{N-1} \cdot \sum_{k=1}^N (c^{(k)})^2}. \quad (8)$$

Computation time required for Monte-Carlo simulation. The relative error of the above statistical estimate depends only on N (as $\approx 1/\sqrt{N}$), and not on the number of variables n . Therefore, the number N_f of calls to f that is needed to achieve a given accuracy does not depend on the number of variables at all.

The error of the above algorithm is asymptotically normally distributed, with a standard deviation $\sigma_e \sim \sigma/\sqrt{2N}$. Thus, if we use a “two sigma” bound, we conclude that with probability 95%, this algorithm leads to an estimate for σ which differs from the actual value of σ by $\leq 2\sigma_e = 2\sigma/\sqrt{2N}$.

This is an error with which we estimate the error of indirect measurement; we do not need too much accuracy in this estimation, because, e.g., in real life, we say that an error is $\pm 10\%$ or $\pm 20\%$, but *not* that the error is, say, $\pm 11.8\%$. Therefore, in estimating the error of indirect measurements, it is sufficient to estimate the characteristics of this error with a relative accuracy of, say, 20%.

For the above “two sigma” estimate, this means that we need to select the smallest N for which $2\sigma_e = 2\sigma/\sqrt{2N} \leq 0.2 \cdot \sigma$, i.e., to select $N_f = N = 50$.

In many practical situations, it is sufficient to have a standard deviation of 20% (i.e., to have a “two sigma” guarantee of 40%). In this case, we need only $N = 13$ calls to f .

On the other hand, if we want to guarantee 20% accuracy in 99.9% cases, which correspond to “three sigma”, we must use N for which $3\sigma_e = 3 \cdot \sigma / \sqrt{2N} \leq 0.2 \cdot \sigma$, i.e., we must select $N_f = N = 113$, etc.

For $n \approx 10^3$, all these values of N_f are much smaller than $N_f = n$ required for numerical differentiation.

So, if we have to choose between the (deterministic) numerical differentiation and the randomized Monte-Carlo algorithm, we must select:

- a deterministic algorithm when the number of variables n satisfies the inequality $n \leq N_0$ (where $N_0 \approx 50$), and
- a randomized method if $n \geq N_0$.

Additional advantage: parallelization. In Monte-Carlo algorithm, we need 50 calls to f . If each call requires a minute, the resulting time takes about an hour, which may be too long for on-line results. Fortunately, different calls to the function f are independent on each other, so we can run all the simulations in parallel.

The more processors we have, the less time the resulting computation will take. If we have as many processors as the required number of calls, then the time needed to estimate the error of indirect measurement becomes equal to the time of a single call, i.e., to the time necessary to compute the result \tilde{y} of this indirect measurement. Thus, if we have enough processors working in parallel, we can compute the result of the indirect measurement *and* estimate its error during the same time that it normally takes just to compute the result.

In particular, if the result \tilde{y} of indirect measurement can be computed in real time, we can estimate the error of this result in real time as well.

5 Solution for Interval and Fuzzy Setting: New Method Based on Cauchy Distribution

Can we use a similar idea in the interval setting? Since Monte-Carlo simulation speeds up computations, it is desirable to use a similar technique in interval setting as well.

There is a problem here. In the interval setting, we do not know the exact distribution, we may have different probability distributions – as long as they are located within the corresponding intervals. If we only use one of these distributions for simulations, there is no guarantee that the results will be valid for other distributions as well.

In principle, we could repeat simulations for several different distributions, but this repetition would drastically increase the simulation time and thus, eliminate the advantages of simulation as opposed to numerical differentiation.

Yes, we can. Luckily, there is a mathematical trick that enables us to use Monte-Carlo simulation in interval setting as well. This trick is based on using *Cauchy distribution* – i.e., probability distributions with the probability density

$$\rho(z) = \frac{\Delta}{\pi \cdot (z^2 + \Delta^2)}; \quad (9)$$

the value Δ is called the *scale parameter* of this distribution, or simply a *parameter*, for short.

Cauchy distribution has the following property that we will use: if z_1, \dots, z_n are independent random variables, and each of z_i is distributed according to the Cauchy law with parameter Δ_i , then their linear combination $z = c_1 \cdot z_1 + \dots + c_n \cdot z_n$ is also distributed according to a Cauchy law, with a scale parameter $\Delta = |c_1| \cdot \Delta_1 + \dots + |c_n| \cdot \Delta_n$.

Therefore, if we take random variables δ_i which are Cauchy distributed with parameters Δ_i , then the value

$$c = f(\tilde{x}_1 + \delta_1, \dots, \tilde{x}_n + \delta_n) - f(\tilde{x}_1, \dots, \tilde{x}_n) = c_1 \cdot \delta_1 + \dots + c_n \cdot \delta_n \quad (10)$$

is Cauchy distributed with the desired parameter (4). So, repeating this experiment N times, we get N values $c^{(1)}, \dots, c^{(N)}$ which are Cauchy distributed with the unknown parameter, and from them we can estimate Δ .

The bigger N , the better estimates we get.

There are two questions to be solved:

- how to simulate the Cauchy distribution;
- how to estimate the parameter Δ of this distribution from a finite sample.

Simulation can be based on the functional transformation of uniformly distributed sample values:

$$\delta_i = \Delta_i \cdot \tan(\pi \cdot (r_i - 0.5)), \quad (11)$$

where r_i is uniformly distributed on the interval $[0, 1]$.

In order to estimate σ , we can apply the Maximum Likelihood Method $\rho(d^1) \cdot \rho(d^2) \cdot \dots \cdot \rho(d^n) \rightarrow \max$, where $\rho(z)$ is a Cauchy distribution density with the unknown Δ . When we substitute the above-given formula for $\rho(z)$ and equate the derivative of the product with respect to Δ to 0 (since it is a maximum), we get an equation

$$\frac{1}{1 + \left(\frac{c^{(1)}}{\Delta}\right)^2} + \dots + \frac{1}{1 + \left(\frac{c^{(N)}}{\Delta}\right)^2} = \frac{N}{2}. \quad (12)$$

The left-hand side of (12) is an increasing function that is equal to 0 ($< N/2$) for $\Delta = 0$ and $> N/2$ for $\Delta = \max |c^{(k)}|$; therefore the solution to the equation (12) can be found by applying a bisection method to the interval $[0, \max |c^{(k)}|]$.

It is important to mention that we assumed that the function f is reasonably linear within the box

$$[\tilde{x}_1 - \Delta_1, \tilde{x}_1 + \Delta_1] \times \dots \times [\tilde{x}_n - \Delta_n, \tilde{x}_n + \Delta_n]. \quad (13)$$

However, the simulated values δ_i may be outside the box. When we get such values, we do not use the function f for them, we use a normalized function that is equal to f within the box, and that is extended linearly for all other values (we will see, in the description of an algorithm, how this is done).

As a result, we arrive at the following algorithm (described, for a somewhat different problem, in [32, 33, 36, 59]):

Algorithm.

- Apply f to the results of direct measurements: $\tilde{y} := f(\tilde{x}_1, \dots, \tilde{x}_n)$;
- For $k = 1, 2, \dots, N$, repeat the following:
 - use the standard random number generator to compute n numbers $r_i^{(k)}$, $i = 1, 2, \dots, n$, that are uniformly distributed on the interval $[0, 1]$;
 - compute Cauchy distributed values $c_i^{(k)} := \tan(\pi \cdot (r_i^{(k)} - 0.5))$;
 - compute the largest value of $|c_i^{(k)}|$ so that we will be able to normalize the simulated measurement errors and apply f to the values that are within the box of possible values: $K := \max_i |c_i^{(k)}|$;
 - compute the simulated measurement errors $\delta_i^{(k)} := \Delta_i \cdot c_i^{(k)} / K$;
 - compute the simulated measurement results $x_i^{(k)} := \tilde{x}_i + \delta_i^{(k)}$;
 - apply the program f to the simulated measurement results and compute the simulated error of the indirect measurement:

$$c^{(k)} := K \cdot \left(f(x_1^{(k)}, \dots, x_n^{(k)}) - \tilde{y} \right);$$

- Compute Δ by applying the bisection method to solve the equation (12).

When is this randomized algorithm better than deterministic numerical differentiation? To determine the parameter Δ , we use the maximum likelihood method. It is known that the error of this method is asymptotically normally distributed, with 0 average and standard deviation $1/\sqrt{N \cdot I}$, where I is Fisher's information:

$$I = \int_{-\infty}^{\infty} \frac{1}{\rho} \cdot \left(\frac{\partial \rho}{\partial \Delta} \right)^2 dz.$$

For Cauchy probability density $\rho(z)$, we have $I = 1/(2\Delta^2)$, so the error of the above randomized algorithm is asymptotically normally distributed, with a standard deviation $\sigma_e \sim \Delta \cdot \sqrt{2/N}$. Thus, if we use a “two sigma” bound, we conclude that with probability 95%, this algorithm leads to an estimate for Δ which differs from the actual value of Δ by $\leq 2\sigma_e = 2\Delta \cdot \sqrt{2/N}$. So, if

we want to achieve a 20% accuracy in the error estimation, we must use the smallest N for which $2\sigma_e = 2\Delta \cdot \sqrt{2/N} \leq 0.2 \cdot \Delta$, i.e., to select $N_f = N = 200$.

When it is sufficient to have a standard deviation of 20% (i.e., to have a “two sigma” guarantee of 40%), we need only $N = 50$ calls to f . For $n \approx 10^3$, both values N_f are much smaller than $N_f = n$ required for numerical differentiation.

So, if we have to choose between the (deterministic) numerical differentiation and the randomized Monte-Carlo algorithm, we must select:

- a deterministic algorithm when the number of variables n satisfies the inequality $n \leq N_0$ (where $N_0 \approx 200$), and
- a randomized algorithm if $n \geq N_0$.

Comment. If we use fewer than N_0 simulations, then we still get an approximate value of the range, but with worse accuracy – and the accuracy can be easily computed by using the above formulas.

This algorithm is naturally parallelizable. Similarly to the Monte-Carlo algorithm for statistical setting, we can run all N simulations in parallel and thus, speed up the computations.

Conclusion. When we know the code for f , then we can use AD and compute Δ and σ in time $O(T + n)$. If the owner of the program f only allows to use it as a black box, then we cannot use AD any more. In principle, we can compute each of n derivatives $\partial f / \partial x_i$ by numerical differentiation, but this would require computation time $T \cdot n \gg T + n$.

For probabilistic uncertainty, one can use Monte-Carlo simulations and compute σ in time $O(T) \ll T \cdot n$. We have shown that for interval uncertainty, we can also compute Δ in time $O(T)$ by using an artificial Monte-Carlo simulations in which each Δx_i is Cauchy distributed with parameter Δ_i – then simulated Δy is Cauchy distributed with the desired parameter Δ .

Remark: the problem of non-linearity. In the above text, we assumed that the intervals \mathbf{x}_i are narrow. In this case, terms quadratic in Δx_i are negligible, and so, we can safely assume that the desired function $f(x_1, \dots, x_n)$ is linear on the box

$$\mathbf{x}_1 \times \dots \times \mathbf{x}_n.$$

In practice, some intervals \mathbf{x}_i may be wide, so even when restricted to the box, the function $f(x_1, \dots, x_n)$ is non-linear. What can we do in this case?

Usually, experts (e.g., designers of the corresponding technical system) know for which variables x_i , the dependence is non-linear. For each of these variables, we can *bisect* the corresponding interval $[\underline{x}_i, \bar{x}_i]$ into two smaller subintervals – for which the dependence is approximately linear. Then, we estimate the range of the function f separately on each of the resulting sub-boxes, and take the union of these two ranges as the range over the entire box.

If one bisection is not enough and the dependence of f on x_i is non-linear over one or several subboxes, we can bisect these boxes again, etc.

This bisection idea has been successfully used in interval computations; see, e.g., [24, 26].

6 Summary

In different knowledge domains, there is a large amount of data stored in different locations; algorithms for processing this data are also implemented at different locations. Web services – and, more generally, cyberinfrastructure – provide the users with an efficient way to submit requests without worrying about the geographic locations of different computational resources (databases and programs) – and avoid centralization with its excessive workloads [21]. Web services enable the user to receive the desired data x_1, \dots, x_n and the results $y = f(x_1, \dots, x_n)$ of processing this data.

The data x_i usually come from measurements or from experts. Measurements are never 100% accurate; as a result, the measured values \tilde{x}_i are, in general, somewhat different from the actual (unknown) values x_i of the corresponding quantities. Experts can also only provide us with approximate values of the desired quantities.

As a result of this measurement or expert uncertainty, the result $\tilde{y} = f(\tilde{x}_1, \dots, \tilde{x}_n)$ of data processing is, in general, different from the actual value $y = f(x_1, \dots, x_n)$ of the desired quantity. It is desirable to gauge this difference.

Traditional methods for estimating the resulting uncertainty in y are based on the assumption that we know the code of the function f . In cyberinfrastructure, owners of the program f may not want to disclose its code; instead, they may only allow to use f as a black box. In this case, traditional techniques are not applicable.

There exist techniques for processing uncertainty under such a “black-box” situation, but these techniques require much longer computation time. In this chapter, we describe new Monte-Carlo-type techniques that process uncertainty in such privacy-protecting black-box situations and that require the same amount of computation time as the traditional non-privacy-protecting techniques.

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 thankful to the anonymous referees for valuable suggestions.

References

1. Aguiar MS, Dimuro GP, Costa ACR, Silva RKS, Costa FA, Kreinovich V (2004) The multi-layered interval categorizer tessellation-based model, In: Iochpe C, Câmara G (eds), IFIP WG2.6 Proceedings of the 6th Brazilian Symposium on Geoinformatics Geoinfo'2004, Campos do Jordão, Brazil, November 22–24, 2004, pp. 437–454. ISBN 3901882200
2. Aldouri R, Keller GR, Gates A, Rasillo J, Salayandia L, Kreinovich V, Seeley J, Taylor P, Holloway S (2004) GEON: Geophysical data add the 3rd dimension in geospatial studies. In: Proceedings of the ESRI International User Conference 2004, San Diego, California, August 9–13, 2004, Paper 1898
3. Averill MG, Miller KC, Keller GR, Kreinovich V, Araiza R, Starks SA (2005) Using expert knowledge in solving the seismic inverse problem, In: Proceedings of the 24nd International Conference of the North American Fuzzy Information Processing Society NAFIPS'2005, Ann Arbor, Michigan, June 22–25, 2005, pp. 310–314
4. Averill MG, Miller KC, Keller GR, Kreinovich V, Araiza R, Starks SA (2007) Using Expert Knowledge in Solving the Seismic Inverse Problem. *International Journal of Approximate Reasoning*, to appear
5. Bardossy G, Fodor J (2004) Evaluation of uncertainties and risks in geology. Springer Verlag, Berlin
6. Berz M, Bischof C, Corliss G, Griewank A (1996), Computational differentiation: techniques, applications, and tools. SIAM, Philadelphia
7. Bojadziev G, Bojadziev M (1995) Fuzzy sets, fuzzy logic, applications. World Scientific, Singapore
8. Ceberio M, Ferson S, Kreinovich V, Chopra S, Xiang G, Murguia A, Santillan J (2006) How to take into account dependence between the inputs: from interval computations to constraint-related set computations, with potential applications to nuclear safety, bio- and geosciences. In: Proceedings of the Second International Workshop on Reliable Engineering Computing, Savannah, Georgia, February 22–24, 2006, pp. 127–154
9. Ceberio M, Kreinovich V, Chopra S, Ludäscher B (2005) Taylor model-type techniques for handling uncertainty in expert systems, with potential applications to geoinformatics. In: Proceedings of the 17th World Congress of the International Association for Mathematics and Computers in Simulation IMACS'2005, Paris, France, July 11–15, 2005
10. Dalenius T (1986) Finding a needle in a haystack – or identifying anonymous census record. *Journal of Official Statistics* 2(3):329–336
11. Demicco R, Klir G, eds (2003) Fuzzy logic in geology. Academic Press
12. Denning DERD (1982) Cryptography and data security. Addison-Wesley, Reading, Massachusetts
13. Doser DI, Crain KD, Baker MR, Kreinovich V, Gerstenberger MC (1998) Estimating uncertainties for geophysical tomography. *Reliable Computing* 4(3):241–268
14. Dubois D, Prade H (1978) Operations on fuzzy numbers. *International Journal of Systems Science* 9:613–626
15. Dubois D, Prade H (1979) Fuzzy real algebra: some results. *Fuzzy Sets and Systems* 2:327–348
16. Dubois D, Prade H (1980) Fuzzy sets and systems: theory and applications. Academic Press, New York, London

17. Duncan G, Lambert D (1987) The risk of disclosure for microdata. In: *Proc. of the Bureau of the Census Third Annual Research Conference*, Bureau of the Census, Washington, DC, 263–274
18. Duncan G, Mukherjee S (1991) Microdata disclosure limitation in statistical databases: query size and random sample query control. In: Prof. 1991 IEEE Symposium on Research in Security and Privacy, Oakland, CA, May 20–22, 1991
19. Fellegi I (1972) On the question of statistical confidentiality. *Journal of the American Statistical Association* 7–18
20. Fuller WA (1987) *Measurement error models*. J. Wiley & Sons, New York
21. Gates A, Kreinovich V, Longpré L, Pinheiro da Silva P, Keller GR (2006) Towards secure cyberinfrastructure for sharing border information. In: *Proceedings of the Lineae Terrarum: International Border Conference*, El Paso, Las Cruces, and Cd. Juárez, March 27–30, 2006
22. Griewank A (2000) *Evaluating derivatives: principles and techniques of algorithmic differentiation*. SIAM, Philadelphia
23. Hole JA (1992) Nonlinear high-resolution three-dimensional seismic travel time tomography. *J. Geophysical Research* 97(B5):6553–6562.
24. Jaulin L, Kieffer M, Didrit O, Walter E (2001) *Applied interval analysis*. Springer Verlag, London
25. Kauffman A, Gupta MM (1985) *Introduction to fuzzy arithmetic: theory and applications*. Van Nostrand, New York
26. Kearfott RB, Kreinovich V, eds (1996) *Applications of interval computations*. Kluwer, Dordrecht
27. Keller GR, Hildenbrand TG, Kucks R, Webring M, Briesacher A, Rujawitz K, Hittleman AM, Roman DJ, Winester D, Aldouri R, Seeley J, Rasillo J, Torres T, Hinze WJ, Gates A, Kreinovich V, Salayandia L (2006) A community effort to construct a gravity database for the United States and an associated Web portal. In: Sinha AK (ed), *Geoinformatics: Data to Knowledge*, Geological Society of America Publ., Boulder, Colorado, pp. 21–34
28. Kim J (1986) A method for limiting disclosure of microdata based on random noise and transformation. In: *Proceedings of the Section on Survey Research Methods of the American Statistical Association* 370–374
29. Kirkendall N et al. (1994) *Report on Statistical Disclosure Limitations Methodology*. Office of Management and Budget, Washington, DC, Statistical Policy Working Paper No. 22
30. Klir G, Yuan B (1995) *Fuzzy sets and fuzzy logic*. Prentice Hall, New Jersey
31. Kreinovich V, Beck J, Ferregut C, Sanchez A, Keller GR, Averill M, Starks SA (2007) Monte-Carlo-type techniques for processing interval uncertainty, and their potential engineering applications. *Reliable Computing* 13(1):25–69.
32. Kreinovich V, Bernat A, Villa E, Mariscal Y (1991) Parallel computers estimate errors caused by imprecise data. *Interval Computations* (2):21–46
33. Kreinovich V, Ferson S (2004) A new cauchy-based black-box technique for uncertainty in risk analysis. *Reliability Engineering and Systems Safety* 85(1–3):267–279
34. Kreinovich V, Lakeyev A, Rohn J, Kahl P (1998) *Computational complexity and feasibility of data processing and interval computations*. Kluwer, Dordrecht
35. Kreinovich V, Longpré L, Starks SA, Xiang G, Beck J, Kandathi K, Nayak A, Ferson S, Hajagos J (2007) Interval versions of statistical techniques, with

- applications to environmental analysis, bioinformatics, and privacy in statistical databases. *Journal of Computational and Applied Mathematics* 199(2):418–423
36. Kreinovich V, Pavlovich MI (1985) Error estimate of the result of indirect measurements by using a calculational experiment. *Measurement Techniques* 28(3):201–205
 37. Longpré L, Kreinovich V, Freudenthal E, Ceberio M, Modave F, Baijal N, Chen W, Chirayath V, Xiang G, Vargas JI (2005) Privacy: protecting, processing, and measuring loss. In: Abstracts of the 2005 South Central Information Security Symposium SCISS'05, Austin, Texas, April 30, 2005, p. 2
 38. Moore RE, Lodwick WA (2003) Interval analysis and fuzzy set theory. *Fuzzy Sets and Systems* 135(1):5–9
 39. Morgenstern M (1987) Security and inference in multilevel database and knowledge base systems. In: *Proc. of the ACM SIGMOD Conference* 357–373
 40. Nguyen HT, Kreinovich V (1996) Nested intervals and sets: concepts, relations to fuzzy sets, and applications, In: Kearfott RB, Kreinovich V (eds) *Applications of interval computations*. Kluwer, Dordrecht, pp. 245–290.
 41. Nguyen HT, Walker EA (2005) *First course in fuzzy logic*. CRC Press, Boca Raton, Florida
 42. Office of Technology Assessment (1993) *Protecting privacy in computerized medical information*. US Government Printing Office, Washington, DC
 43. Palley M, Siminoff J (1986) Regression methodology based disclosure of a statistical database. In: *Proceedings of the Section on Survey Research Methods of the American Statistical Association* 382–387
 44. Parker RL (1994) *Geophysical inverse theory*. Princeton University Press, Princeton, New Jersey
 45. Platon E, Tupelly K, Kreinovich V, Starks SA, Villaverde K (2005) Exact bounds for interval and fuzzy functions under monotonicity constraints, with potential applications to biostratigraphy. In: *Proceedings of the 2005 IEEE International Conference on Fuzzy Systems FUZZ-IEEE'2005*, Reno, Nevada, May 22–25, 2005, pp. 891–896
 46. Rabinovich S (2005) *Measurement errors and uncertainties: theory and practice*. American Institute of Physics, New York
 47. Schiek CG, Araiza R, Hurtado JM, Velasco AA, Kreinovich V, Sinyansky V (2006) Images with uncertainty: efficient algorithms for shift, rotation, scaling, and registration, and their applications to geosciences. In: Nachtigael M, Van der Weken D, Kerre EE (eds), *Soft computing in image processing: recent advances*, Springer Verlag
 48. Sinha AK (2006) *Geoinformatics: data to knowledge*. Geological Society of America Publ., Boulder, Colorado
 49. Su T, Ozsoyoglu G (1991) Controlling FD and MVD inference in multilevel relational database systems. *IEEE Transactions on Knowledge and Data Engineering* 3:474–485
 50. Sweeney L (1996) Replacing personally-identifying information in medical records, the scrub system. *Journal of the American Medical Informatics Association* 333–337
 51. Sweeney L (1997) Weaving technology and policy together to maintain confidentiality. *Journal of Law, Medicine and Ethics* 25:98–110
 52. Sweeney L (1997) Guaranteeing anonymity when sharing medical data, the datafly system. *Journal of the American Medical Informatics Association* 51–55

53. Sweeney L (1997) Computational disclosure control for medical microdata. In: Proceedings of the Record Linkage Workshop, Bureau of the Census, Washington, DC
54. Sweeney L (1998) Commentary: researchers need not rely on consent or not. *New England Journal of Medicine* 338(15)
55. Sweeney L (1998) Towards the optimal suppression of details when disclosing medical data, the use of sub-combination analysis. In: Proceedings of MED-INFO'98, International Medical Informatics Association, Seoul, Korea, North-Holland, p. 1157
56. Sweeney L (1998) Three computational systems for disclosing medical data in the year 1999. In: Proceedings of MEDINFO'98, International Medical Informatics Association, Seoul, Korea, North-Holland pp. 1124–1129
57. Sweeney L (1998) Datafly: a system for providing anonymity in medical data. In: Lin TY, Qian S (eds.) *Database security XI: status and Prospects*. Elsevier, Amsterdam
58. Torres R, Keller GR, Kreinovich V, Longpré L, Starks SA (2004) Eliminating duplicates under interval and fuzzy uncertainty: an asymptotically optimal algorithm and its geospatial applications. *Reliable Computing* 10(5):401–422
59. Trejo R, Kreinovich V (2001) Error estimations for indirect measurements: randomized vs. deterministic algorithms for 'black-box' programs. In: Rajasekaran S, Pardalos P, Reif J, and Rolim J (eds), *Handbook on randomized computing*. Kluwer, Dordrecht, 673–729
60. Vavasis SA (1991) *Nonlinear optimization: complexity issues*, Oxford University Press, New York
61. Willenborg L, De Waal T (1996) *Statistical disclosure control in practice*. Springer Verlag, New York
62. Wen Q, Gates AQ, Beck J, Kreinovich V, Keller JR (2001) Towards automatic detection of erroneous measurement results in a gravity database. In: Proceedings of the 2001 IEEE Systems, Man, and Cybernetics Conference, Tucson, Arizona, October 7–10, 2001, pp. 2170–2175
63. Xie H, Hicks N, Keller GR, Huang H, Kreinovich V (2003) An IDL/ENVI implementation of the FFT based algorithm for automatic image registration. *Computers and Geosciences* 29(8):1045–1055
64. Zelt CA, Barton PJ (1998) Three-dimensional seismic refraction tomography: A comparison of two methods applied to data from the Faeroe Basin. *J. Geophysical Research* 103(B4):7187–7210.