

## **Representation and Problem Solving with Distribution Envelope Determination (DEnv)**

Daniel Berleant and Jianzhong Zhang  
Department of Electrical and Computer Engineering  
2215 Coover Hall  
Iowa State University  
Ames, Iowa 50011  
berleant@iastate.edu

### **Abstract**

Distribution Envelope Determination (DEnv) is a method for computing the CDFs of random variables whose samples are a function of samples of other random variable(s), termed inputs. DEnv computes envelopes around these CDFs when there is uncertainty about the precise form of the probability distribution describing any input. For example, inputs whose distribution functions have means and variances known only to within intervals can be handled. More generally, inputs can be handled if the set of all plausible cumulative distributions describing each input can be enclosed between left and right envelopes. Results will typically be in the form of envelopes when inputs are envelopes, when the dependency relationship of the inputs is unspecified, or both. For example in the case of specific input distribution functions with unspecified dependency relationships, each of the infinite number of possible dependency relationships would imply some specific output distribution, and the set of all such output distributions can be bounded with envelopes. The DEnv algorithm is a way to obtain the bounding envelopes. DEnv is implemented in a tool which is used to solve problems from a benchmark set.

**Keywords.** DEnv, p-boxes, aleatory uncertainty, epistemic uncertainty, second order uncertainty, uncertainty quantification, 2<sup>nd</sup> order uncertainty, reducible uncertainty, imprecise probabilities, challenge problems, envelopes, derived distributions, Statool.

### **1 Introduction**

The DEnv (Distribution Envelope Determination) algorithm is a method for computing distributions whose samples are some function of the samples of other input distributions, even under non-traditional conditions of severely limited knowledge about the inputs.

Under traditional conditions of known dependency relationships among precisely defined input distributions, solutions based around Monte Carlo simulation have an extensive literature, although MC gives results that are potentially problematic (Ferson 1996 [9]) and whose interpretation can be complicated by random variation especially in tails and other unlikely regions of system behavior. A large literature also addresses analytical solutions, which tend to require certain well-defined classes of distributions as inputs (Springer 1979 [28] is fairly comprehensive up to its time of writing). When the input random variables to be combined are independent in the traditional sense that the probability of a joint event is the product of the probabilities of its constituent events (often termed stochastic [7] or statistical independence), solutions based on numerical convolution are well known (Ingram et al. 1968 [16], Colombo and Jaarsma 1980 [6], Kaplan 1981 [17], Moore 1984 [20]). Lodwick's (2003 [19]) method is applied to multivariate examples with repeating variables and stated to be usable when variables are non-independent, or when their dependencies are unspecified, which is among the following problem characteristics that pose a challenge to traditional approaches.

- (1) Sample values of one of the random variables may be described by a distribution, while sample values of another may be known only to within an interval.
- (2) The input random variables may not be independent, and their dependency relationship may be unknown or only partly known. We will use the term *unknown dependency* to describe this situation. The term “unknown interaction” has also been used (Couso et al. 1999).
- (3) There may be insufficient information available to assign a specific distribution to an input random variable.

The problem of combining a distribution with an interval, (1) above, was addressed by Berleant (1993 [1]). When input dependencies are unknown, (2) above, the result random variable cannot in general be described with a single distribution, because each possible dependency relationship between the inputs leads to its own result distribution. Frank et al. 1987 [14] discuss the distribution of sums and products of samples of other distributions under this condition. Envelopes, also called p(robability)-bounds or p(robability)-boxes (Ferson et al. [10]) can be found which surround the family of possible result distributions. If these envelopes around the results are to be used in turn as inputs to produce further results, the algorithm for obtaining the further results must be able to use envelopes as inputs. This is also the problem of (3) above. We review solutions to (2), and then (3), next.

One approach to manipulating envelopes and distributions with unknown dependency relationships is based on the Probabilistic Arithmetic of Williamson and Downs (1990 [30]), which in turn is built on a foundation of copulas (Nelsen 1999 [21]). Probabilistic Arithmetic is one component of the commercially available RiskCalc (Ferson 2002 [8]) software. An approach based on sets of probability measures was applied to problems from a benchmark set (Oberkampf et al., this issue [23]) by Fetz and Oberguggenberger (this issue [13]), as was a Monte-Carlo based approach (Red-Horse this issue [24]). Ferson and Hajagos (this issue [11]) also address the problems using the just-mentioned Probabilistic Arithmetic. Tonon (this issue [29]) addresses Problem B using random set theory. Further solutions and insights were also presented by others at a recent workshop (Sandia 2002 [26]). In related work Neumaier [22] recently described clouds, a concept capable of expressing and manipulating families of CDFs bounded by left and right envelopes. The approach described in this paper is Distribution Envelope Determination (DEnv), which relies on safely discretized distributions and linear programming. It has been reported on a theoretical basis (Berleant and Goodman-Strauss 1998 [2]) and implemented in a tool (Berleant et al. 2003 [3]). Applications have also been described (Berleant et al. 2002 [4]; Sheblé and Berleant 2002 [27]).

Equivalence properties of DEnv, Probabilistic Arithmetic, imprecise probabilities, and Dempster-Shafer structures are described by Regan et al. [25]. It appears these approaches are largely equivalent in their ability to construct envelopes around cumulative distributions in the real domain. They are also extendable to fuzzy numbers. Questions about how they compare in terms of computational speed and in ability to express and use inputs that are in non-cumulative form have still not been fully resolved. We feel that DEnv has an advantage in understandability compared to other methods. For example Probabilistic Arithmetic requires an understanding of copulas. Random sets also require specialized knowledge. Although DEnv uses linear programming (LP), knowledge of LP is widespread, and in DEnv may be viewed as a black box.

## 1.1 Concise review of the DEnv Algorithm

This section describes DEnv concisely and abstractly. Section 2 covers DEnv less formally, but at length and in detail in the context of a set of challenge problems [23]. The reader may choose to skip directly to section 2 without loss of continuity and refer back to this section later as needed, may choose to use this section as a foundation, or may take some intermediate course.

DEnv begins with two inputs, probability density functions  $f_x(\cdot)$  and  $f_y(\cdot)$  describing samples  $x$  and  $y$  of random variables  $X$  and  $Y$ . DEnv will characterize the CDF (cumulative distribution function)  $F_z(z)$  of samples  $z=g(x,y)$  of random variable  $Z$ , given function  $g$ . The input PDFs  $f_x(\cdot)$  and  $f_y(\cdot)$  are discretized by partitioning the support (i.e. the domain over which a PDF is non-zero) of each, yielding intervals  $\mathbf{x}_i$ ,  $i=1\dots I$ , and  $\mathbf{y}_j$ ,  $j=1\dots J$ . Each  $\mathbf{x}_i$  is assigned a probability mass

$$p_{\mathbf{x}_i} = p(x \in \mathbf{x}_i) = \int_{x_0=\underline{\mathbf{x}_i}}^{\overline{\mathbf{x}_i}} f_x(x_0) dx_0, \text{ where interval-valued symbols are shown in bold, and}$$

interval  $\mathbf{x}_i$  has lower bound  $\underline{\mathbf{x}_i}$  and upper bound  $\overline{\mathbf{x}_i}$ . Similarly, each  $\mathbf{y}_j$  is assigned a probability

$$\text{mass } p_{\mathbf{y}_j} = p(y \in \mathbf{y}_j) = \int_{y_0=\underline{\mathbf{y}_j}}^{\overline{\mathbf{y}_j}} f_y(y_0) dy_0. \text{ The } \mathbf{x}_i\text{'s and } \mathbf{y}_j\text{'s and their probabilities form the}$$

marginals of a discretized joint distribution called a *joint distribution tableau* (Table 1), the interior cells of which each contain two items. One is a probability mass

$$p_{ij} = p(x \in \mathbf{x}_i \wedge y \in \mathbf{y}_j). \text{ If the value of } x \text{ gives no information about the value of } y, \text{ and vice}$$

versa, then  $x$  and  $y$  are independent and  $p_{ij} = p(x \in \mathbf{x}_i) \cdot p(y \in \mathbf{y}_j) = p_{\mathbf{x}_i} \cdot p_{\mathbf{y}_j}$ , where  $p_{\mathbf{x}_i}$  is

defined as  $p(x \in \mathbf{x}_i)$  and  $p_{\mathbf{y}_j}$  as  $p(y \in \mathbf{y}_j)$ . The second item is an interval that bounds the

values  $z=g(x,y)$  may have, given  $x \in \mathbf{x}_i \wedge y \in \mathbf{y}_j$ . In other words,  $\mathbf{z}_{ij}=\mathbf{g}(\mathbf{x}_i, \mathbf{y}_j)$ .

$x \rightarrow$				
$y \downarrow$	$z=g(x,y) \searrow$	$\cdots$	$\mathbf{x}_i$	$\cdots$
			$p_{\mathbf{x}_i} = p(x \in \mathbf{x}_i)$	
$\vdots$		$\ddots$	$\vdots$	$\ddots$
$\mathbf{y}_j$		$\cdots$	$\mathbf{z}_{ij}=\mathbf{g}(\mathbf{x}_i, \mathbf{y}_j)$	$\cdots$
$p_{\mathbf{y}_j} = p(y \in \mathbf{y}_j)$			$p_{ij} = p(x \in \mathbf{x}_i \wedge y \in \mathbf{y}_j)$	
$\vdots$		$\ddots$	$\vdots$	$\ddots$

**Table 1. General form of a joint distribution tableau.**

To better characterize the CDF  $F_z(\cdot)$ , we next convert the set of interior cells of the joint distribution tableau into cumulative form. Because the distribution of each probability mass  $p_{ij}$  over its interval  $\mathbf{z}_{ij}$  is not defined by the tableau, values of  $F_z(\cdot)$  cannot be computed precisely. However they can be bounded, so DEnv does this by computing the interval-valued function  $\mathbf{F}_z(\cdot)$  as

$$\underline{\mathbf{F}}_z(z_0) = \sum_{i,j | \mathbf{z}_{ij} \leq z_0} p_{ij} \text{ and } \overline{\mathbf{F}}_z(z_0) = \sum_{i,j | \mathbf{z}_{ij} \leq z_0} p_{ij}, \quad (1)$$

resulting in right and left envelopes respectively for  $\mathbf{F}_z(\cdot)$ .

An additional complication occurs if the dependency relationship between  $x$  and  $y$  is unknown, because then the values of the  $p_{ij}$ 's are underdetermined and so Equations (1) cannot be evaluated.

However, the  $p_{ij}$ 's in a column of a joint distribution tableau must sum to  $p_{x_i}$  and the  $p_{ij}$ 's in a row must sum to  $p_{y_j}$  giving three sets of constraints:  $p_{x_i} = \sum_j p_{ij}$ ,  $p_{y_j} = \sum_i p_{ij}$ , and  $p_{ij} \geq 0$ , for  $i=1 \dots I, j=1 \dots J$ . These constraints are all linear, and so may be passed to a linear programming routine. Linear programming takes as input linear constraints on variables, which in this case are the  $p_{ij}$ 's, and an expression in those variables to minimize, for example,  $\underline{F}_z(z_0) = \sum_{i,j | z_{ij} \leq z_0} p_{ij}$  for some given value  $z_0$ . The output would then be the minimum value that  $\underline{F}_z(z_0)$  could have such that the values assigned to the  $p_{ij}$ 's are consistent with the constraints.  $\overline{F}_z(z_0) = \sum_{i,j | z_{ij} \leq z_0} p_{ij}$  is maximized similarly. These envelopes are less restrictive (i.e. farther apart) than when the  $p_{ij}$ 's are determined by an assumption of independence or some other dependency relationship so that linear programming is not needed.

These ideas generalize naturally to  $n$  marginals, which would require an  $n$ -dimensional joint distribution tableau (Section 2.6).

## 2 The Challenge Problems and the DEnv Technique

In this section the DEnv technique is explained in the context of the challenge problems given by Oberkampf et al. [23]. Solutions are presented and explained for all of the challenge problems, which include six scenarios involving computation of  $y=(a+b)^a$  and also a spring problem.

### Problem 1: setting the stage

Problem 1 is to find the range of  $y=(a+b)^a$  given  $a \in [0.1, 1]$  and  $b \in [0, 1]$ . The minimum for  $y$  occurs when  $a=0.37$ , which is not an endpoint of the interval for  $a$ , and when  $b=0$ , leading to the answer  $y \in [0.69, 2]$ . Using only endpoints of input intervals to compute bounds on result intervals can, as in this case, generate misleading results. This is a well-studied issue in interval computing and occurs numerous times in the challenge problems with respect to intervals for  $a$  and in the spring problem.

Challenge Problems 2-6 have, as givens, one or more sources of information about  $a$  and also about  $b$ . The sources often are specified to have equal credibility. The equal credibility stipulation contains significant ambiguity. This has serious implications for the solutions, which are discussed later in Section 3. In the present section we seek solutions for the problems and hence must precisely define them. Therefore we resolve the ambiguity by modeling credibility using probability.

When different information sources have equal credibility we interpret this to mean that the actual but unknown value has the same probability of being binned in the interval of one information source as it does of being binned in the interval of any other. This interpretation allows different information sources to have equal credibility while providing probabilities for intervals that are disjoint, nested, or overlapping (all of which could occur in real situations). Which of those situations occurs in a given problem has little effect on the solution using DEnv, an algorithm which is consistent with standard properties of probability and requires, to be applied, that a problem be modeled using intervals and associated probabilities. (Another typical way of modeling problems for processing by DEnv is to discretize probability distributions, as occurs in the solution to Challenge Problem 6.)

## 2.1 Problem 2: $a \in [0.1, 1]$ with equally credible intervals for $b$

The following facts about this problem are defined by Oberkamp et al. (this issue [23]) and the discussion above:  $y = (a+b)^a$ , with  $a \in [0.1, 1]$  and  $p(b \in \mathbf{b}_1) = p(b \in \mathbf{b}_2) = p(b \in \mathbf{b}_3) = p(b \in \mathbf{b}_4)$  for intervals  $\mathbf{b}_1 \dots \mathbf{b}_4$ , each provided by a different information source all of which are equally credible.

Table 2 shows, for Problem 2a, the givens for  $a$  and  $b$  and the results that follow using interval arithmetic to get intervals describing the consequent range of  $y$ . For the last row, the intervals for  $a$ ,  $b$ , and hence  $y$  are the same as for Problem 1 above and hence require the same attention to interior points of  $a$ .

Intervals given for $b$	Intervals for $y = (a+b)^a$ , given $a \in [0.1, 1]$
$\mathbf{b}_1 = [0.6, 0.8]$ $p = 0.25$	$\mathbf{y}_1 = [0.96, 1.8]$ $p_1 = 0.25$
$\mathbf{b}_2 = [0.4, 0.85]$ $p = 0.25$	$\mathbf{y}_2 = [0.9, 1.85]$ $p_2 = 0.25$
$\mathbf{b}_3 = [0.2, 0.9]$ $p = 0.25$	$\mathbf{y}_3 = [0.81, 1.9]$ $p_3 = 0.25$
$\mathbf{b}_4 = [0, 1]$ $p = 0.25$	$\mathbf{y}_4 = [0.69, 2]$ $p_4 = 0.25$

**Table 2.** The interval given by Problem 2a about the value of  $a$  (top right cell), intervals given about the value of  $b$  (left column), and the implications of those intervals for the value of  $y$  (right column).

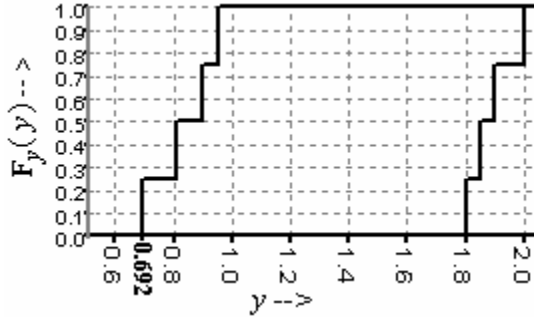
Figure 1 (top) shows  $y$  in Problem 2a graphically. The distribution envelopes shown in the graphs of Figure 1 may be derived straightforwardly from the right-hand column of their corresponding tables as follows.

- 1) *Left envelope.* This envelope,  $\bar{\mathbf{F}}(y)$ , represents the maximum possible cumulation of probability mass for any given value of  $y$ . It is obtained under the extreme assumption that the probability mass associated with each interval is distributed as an impulse at the low bound of the interval. (Such an extreme assumption is permitted because an interval does not imply anything about how its probability mass is distributed beyond that it is distributed within its bounds.) Therefore any interval whose low bound is at or below a value  $y_0$  can contribute up to its full probability mass to the cumulation of  $y$  at  $y_0$ , while other intervals cannot contribute any mass. Thus for each of Problems 2a-2c,  $\bar{\mathbf{F}}(y) = \sum_{k | \underline{\mathbf{y}}_k \leq y} p_k$ , where bold signifies interval-valued symbols, underlining refers to an interval's low bound, and overlining refers to an interval's high bound.

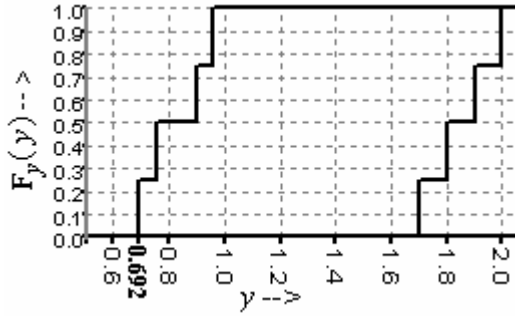
- 2) *Right envelope.* This envelope,  $\underline{\mathbf{F}}(y)$ , represents the minimum cumulation of probability mass for any given value of  $y$ . This is obtained under the extreme assumption that each mass is distributed as an impulse at the high bound of its corresponding interval. Thus,

$$\underline{\mathbf{F}}(y) = \sum_{k | \overline{\mathbf{y}}_k \leq y} p_k.$$

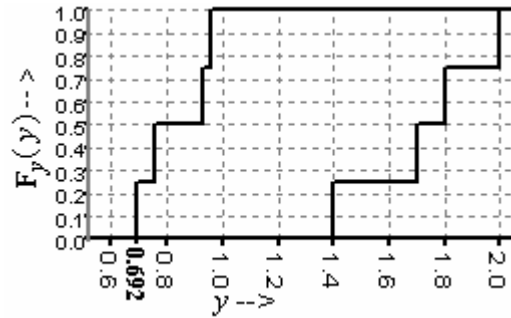
Figure 1 also shows graphs for  $y$  in Problems 2b and 2c, derived as just described. (For Problem 2a, the intervals  $y_k$  were shown in Table 2, while envelopes around the CDF for  $b$  appear later in Figure 15; for Problem 2c, intervals  $y_k$  appear later in Table 4 and envelopes around the CDF of  $b$  appear later in Figure 13, top.) That the intervals for  $b$  in Problem 2a are nested, in 2b are overlapping, and in 2c are completely disjoint [23] does not affect the process of computing intervals and graphs for  $y$ .



(Problem 2a↑)



(Problem 2b↑)



(Problem 2c↑)

Figure 1. Envelopes around the cumulative distribution of the value of  $y$  in Problems 2a, 2b, and 2c implied by the intervals  $y_k$  for each problem.

## 2.2 Problem 3: intervals for $a$ and intervals for $b$

In the preceding problem four sources of information about  $b$  led to computing, then combining, four cases for  $y$ . In this problem, four cases of  $b$  for each of the three cases for  $a$  lead to 12 cases for  $y$ . Given equal probability assignments for each case of  $a$ , and likewise for  $b$ , if  $a$  and  $b$  are assumed independent in the sense used throughout this paper, that a sample of one gives no information about the other, then the 12 cases for  $y$  will each have equal probability. (Fetz and

Oberguggenberger [13] show the implications of different kinds of independence [7] for the challenge problems.) If  $a$  and  $b$  are not independent then cases for  $y$  will typically have different probabilities.

The format of Table 2 can be generalized to express these situations. Figure 2 shows both the table and resulting envelopes for Problem 3a, and the envelopes for 3b (3c will be discussed in detail later). It is convenient to describe the tables using the following terminology.

- The leftmost column and topmost row describe  $a$  and  $b$ , and are called *marginals*.
- Cells of the table with intervals labeled  $y_{ij}$  are called *interior cells*.
- The distribution of the probability mass of a cell over its interval is called the cell's *mini-distribution*.
- The entire table is called a *joint distribution tableau*.

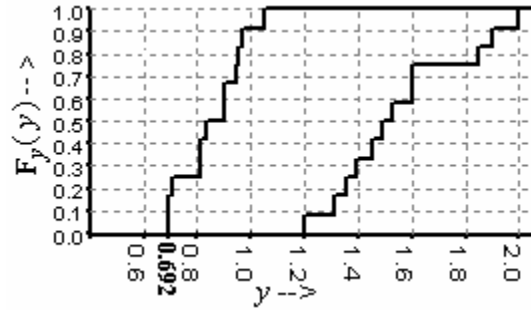
The intervals and probability masses in the marginals discretize  $a$  and  $b$ . The marginal intervals and the interior cell probability masses together discretize the joint distribution of  $a$  and  $b$ . The interior cell probability masses and their intervals give a discretization of the distribution of  $y=(a+b)^a$ . This was shown abstractly in Table 1.

Note that different cells in the same marginal can have overlapping intervals, as in Table 2 and Figure 2 (top). A simple example illustrates the meaning of overlaps. Consider a symmetric probability density function (PDF) with support over  $[0, 4]$ . A very coarse discretization consists of the single interval  $[0, 4]$  with an associated probability mass of 1, because the PDF contains a mass of 1 distributed over the interval  $[0, 4]$ . Since the PDF is symmetric, it may also be discretized as two probability masses, one of 0.5 distributed appropriately over  $[0, 2]$ , and another also of 0.5 and also distributed appropriately over  $(2, 4]$  (assuming the PDF does not have an impulse at exactly 2). Yet, the wider and overlapping intervals  $[0, 3]$  and  $[1, 4]$  can also give a valid discretization of the same PDF. One way to do that is to specify the same probability mass and mini-distribution for  $[0, 3]$  that was just used for  $[0, 2]$  (implying that no mass assigned to  $[0, 3]$  happens to be distributed above 2), and the same mini-distribution over  $[1, 4]$  that was just used for  $(2, 4]$  (implying that no mass assigned to  $[1, 4]$  is distributed at or below 2), resulting in exactly the same mini-distributions as in the case of the  $\{[0, 2], (2, 4]\}$  discretization. As a final example consider a discretization with extreme overlap consisting of two intervals, each with range  $[0, 4]$  and probability mass 0.5. It is certainly possible to distribute one mass within  $[0, 4]$ , then add in the other mass, distributed appropriately within the same interval, to get the original PDF.

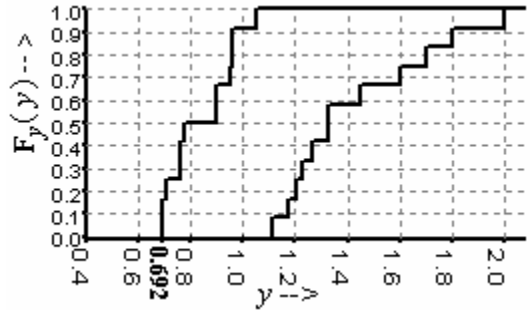
Just as a marginal of a joint distribution tableau discretizes the distribution of an input random variable, the interior cells of the tableau collectively discretize the distribution of  $y=(a+b)^a$  even though the phenomenon of overlapping result intervals is often present (e.g. Table 2 and Figure 2, top) whether or not there are overlapping intervals in either marginal.

$a \rightarrow$ $b \downarrow y \searrow$	[0.5, 0.7] $p=0.33$	[0.3, 0.8] $p=0.33$	[0.1, 1] $p=0.33$
[0.6, 0.6] $p=0.25$	$y_{11}=[1.0, 1.2]$ $p_{11}=0.083$	$y_{21}=[0.97, 1.3]$ $p_{21}=0.083$	$y_{31}=[0.96, 1.6]$ $p_{31}=0.083$
[0.4, 0.85] $p=0.25$	$y_{12}=[0.95, 1.4]$ $p_{12}=0.083$	$y_{22}=[0.90, 1.5]$ $p_{22}=0.083$	$y_{32}=[0.9, 1.85]$ $p_{32}=0.083$
[0.2, 0.9] $p=0.25$	$y_{13}=[0.84, 1.4]$ $p_{13}=0.083$	$y_{23}=[0.81, 1.5]$ $p_{23}=0.083$	$y_{33}=[0.81, 1.9]$ $p_{33}=0.083$
[0, 1] $p=0.25$	$y_{14}=[0.71, 1.4]$ $p_{14}=0.083$	$y_{24}=[0.69, 1.6]$ $p_{24}=0.083$	$y_{34}=[0.69, 2]$ $p_{34}=0.083$

Joint distribution tableau for Problem 3a (numbers are to 2 significant digits).



Envelopes enclosing the CDF for  $y$  in Problem 3a.



Envelopes enclosing the CDF for  $y$  in Problem 3b.

Figure 2. The joint distribution tableau for Problem 3a (top) was used to generate envelopes for  $y$  (middle) using  $\bar{F}(y_0) = \sum_{i,j|y_{ij} \leq y_0} p_{ij}$  for the left envelope and  $\underline{F}(y_0) = \sum_{i,j|y_{ij} \leq y_0} p_{ij}$  for the right envelope (see Section 1.1). Envelopes for  $y$  in Problem 3b are also shown (bottom).

### 2.2.1 Removing the independence assumption

Figure 2 assumes the marginals are independent in the standard statistical sense that the probability assigned to each interior cell in a joint distribution tableau is the product of its marginal cell probabilities. For example, consider the cell in the lower right corner of the tableau of Figure 2 (top) and the marginal cells for its row and its column. The probability that  $a$  is in [0.1, 1] and is binned in the marginal cell with that interval rather than in any other marginal cell whose interval it might be consistent with, is specified as 0.33. Similarly the probability that  $b$  is in the marginal cell with interval [0, 1] is 0.25. Therefore the probability assigned to the lower right interior cell is  $0.25 \cdot 0.33 = 0.083$ .



All interior cell probabilities were computed similarly. However if the marginals are not independent then the interior cell probabilities are determined by the details of the dependency relationship, whatever it is. A human analyst could for example manually type in interior cell probabilities to express some particular dependency relationship, or software could fill them in based on some formula defining a dependency relationship. Equations (1) and the equations in the caption of Figure 2 still apply. The DEnv technique can also be extended to the case where the dependency between  $a$  and  $b$  is not determined. This case has the following two subcases: (i) the dependency between  $a$  and  $b$  is unknown, discussed next, and (ii) partial knowledge about their dependency exists, discussed after.

**Unknown dependency.** If the dependency relationship is not specified then the interior cell probabilities are not determined. They are however constrained by the marginal cell probabilities. The probability in each marginal cell in the left column is distributed among the interior cells in its row. Also the probability in each marginal cell in the top row is distributed among the interior cells in its column. Thus in Table 3 there are four row constraints and three column constraints.

$a \rightarrow$ $b \downarrow \quad y \Downarrow$	$\mathbf{a}_1=[0.8, 1]$ $p=0.33$	$\mathbf{a}_2=[0.5, 0.7]$ $p=0.33$	$\mathbf{a}_3=[0.1, 0.4]$ $p=0.33$
$\mathbf{b}_1=[0.8, 1]$ $p=0.25$	$[1.5, 2]$ $p_{11}=?$	$[1.1, 1.45]$ $p_{21}=?$	$[0.99, 1.1]$ $p_{31}=?$
$\mathbf{b}_2=[0.5, 0.7]$ $p=0.25$	$[1.2, 1.7]$ $p_{12}=?$	$[1.0, 1.3]$ $p_{22}=?$	$[0.93, 1.0]$ $p_{32}=?$
$\mathbf{b}_3=[0.1, 0.4]$ $p=0.25$	$[0.92, 1.4]$ $p_{13}=?$	$[0.78, 1.1]$ $p_{23}=?$	$[0.76, 0.93]$ $p_{33}=?$
$\mathbf{b}_4=[0, 0.2]$ $p=0.25$	$[0.84, 1.2]$ $p_{14}=?$	$[0.71, 0.93]$ $p_{24}=?$	$[0.69, 0.89]$ $p_{34}=?$

Row constraints	Column constraints
$0.25 = p_{11} + p_{21} + p_{31}$	$0.33 = p_{11} + p_{12} + p_{13} + p_{14}$
$0.25 = p_{12} + p_{22} + p_{32}$	$0.33 = p_{21} + p_{22} + p_{23} + p_{24}$
$0.25 = p_{13} + p_{23} + p_{33}$	$0.33 = p_{31} + p_{32} + p_{33} + p_{34}$
$0.25 = p_{14} + p_{24} + p_{34}$	

**Table 3. Joint distribution tableau for Problem 3c (top), expressing the case of unknown dependency between  $a$  and  $b$ . Thus the interior cell probabilities are underdetermined. They are however partially constrained by the marginal cells, each of which defines a row or column constraint (bottom).**

The restrictions on the interior cell probabilities imposed by the row and column constraints are the key to finding the height of the left or right envelope at a given value of  $y=(a+b)^a$ . The formula for the left envelope is

$$\overline{\mathbf{F}}_y(y_0) = \sup_{p_{ij}, i=1 \dots I, j=1 \dots J | C} \sum_{i, j | y_{ij} \leq y_0} p_{ij} \quad \text{and for the right envelope,}$$

$$\underline{\mathbf{F}}_y(y_0) = \inf_{p_{ij}, i=1 \dots I, j=1 \dots J | C} \sum_{i, j | y_{ij} \leq y_0} p_{ij} \quad (2)$$

where  $C$  refers to the set of row and column constraints that must hold. In other words, the sup operation finds probability mass assignments for the  $p_{ij}$ 's that give the maximum value for the summation that is possible while maintaining consistency with constraint set  $C$ , and the inf operation acts similarly to give the minimum value. These equations are like Equations (1) augmented with sup and inf, and express the variability of the  $p_{ij}$ 's as constrained by the row and

column constraints. The intuitions behind these formulas were reviewed toward the end of Section 1.1, and are detailed along with some other salient points next.

*Left envelope.* The height of the left envelope at some value  $y=y_0$  on the horizontal axis is the maximum possible cumulation of probability mass over the interval  $(-\infty, y_0)$ . This may be obtained as follows.

- (1) Identify all interior cells with interval low bounds at or below  $y_0$ . Each potentially contributes its entire probability mass to the cumulation at  $y_0$ , because its mini-distribution can be specified so as to distribute its entire mass over values at or below  $y_0$ . Other interior cells cannot distribute any of their probability masses to values at or below  $y_0$  no matter what their mini-distributions are, because their intervals' low bounds are above  $y_0$ .

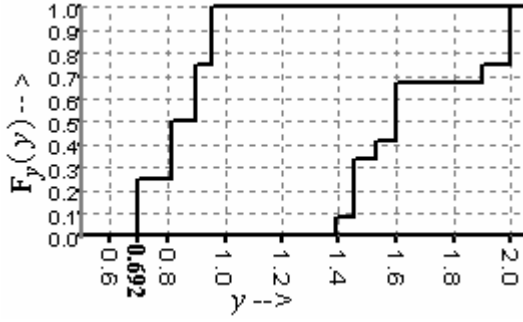
*Example 1.* In Table 3, for  $y=0.95$  the maximum cumulation will involve:

$p_{32}, p_{13}, p_{23}, p_{33}, p_{14}, p_{24}, \& p_{34}.$

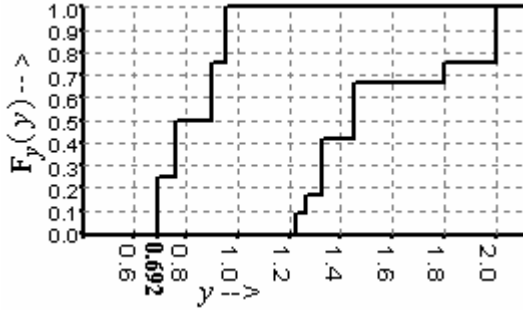
- (2) Maximize the sum of the probability masses of the previously identified interior cells. This may be done by using linear programming as a software subroutine call and passing in as inputs, (i) the constraints defined by the row and column constraints, and (ii) the function to maximize. For *Example 1* above, the function to maximize is therefore  $p_{32}+p_{13}+p_{23}+p_{33}+p_{14}+p_{24}+p_{34}$ . Linear programming finds values for the  $p_{ij}$  which maximize that function while satisfying the constraints. Maximization can actually be done manually by careful inspection, pushing masses around in the tableau, but using a computer to solve this as a linear programming program is more practical.

*Right envelope.* To find the height of the right envelope at  $y=y_0$ , instead of maximizing a sum of interior cell probability masses we minimize, because the right envelope expresses the minimum possible cumulation at each value of  $y$ . The interior cells whose pooled probability mass is to be minimized are those whose intervals have high bounds at or below  $y_0$ , because the full mass of each of those cells must be in the cumulation at  $y_0$ , although from the perspective of minimization we would wish otherwise. It is possible for the mini-distributions of other interior cells to be specified so as to allocate all of their respective masses above  $y_0$ , thereby not contributing to the cumulation. (An alternative to minimizing the mass of this set of interior cells is to maximize the mass of its complement.)

Having explained how to get the height of the left and right envelopes for a given value of  $y$  we must now choose the values of  $y$  at which to do this computation. For the left (right) envelope these values are the low (high) bounds of the interior cell intervals, because it is at these bounds that the envelope heights can change, since the maximization (minimization) process depends on these bounds as described above and in Equations (1). Figure 3 shows the envelopes for Problems 3a and 3b when  $a$  and  $b$  have unknown dependency.



Problem 3a without independence assumption.



Problem 3b without independence assumption.

Figure 3. The envelopes shown here are more widely separated than those of Figure 2, because removing the independence assumption tends to weaken the results.

**Correlation.** Independence is a strong assumption. Simply removing that assumption leaves no information at all about the dependency relationship between  $a$  and  $b$ , significantly weakening results as shown in Figure 3. An intermediate case exists when  $a$  and  $b$  are assumed correlated. Intuitively if  $b$  is likely to be low when  $a$  is low, and high when  $a$  is high, then  $a$  and  $b$  are positively correlated. Alternatively if  $b$  is likely to be high when  $a$  is low, and low when  $a$  is high, then  $a$  and  $b$  are negatively correlated. Consequently in a tableau like that of Table 3 (top), if probability mass is concentrated in interior cells along a diagonal northwest-southeast path the correlation will be high, while if mass is concentrated along the other diagonal correlation will be low. More generally, consider a standard equation for the (Pearson) correlation coefficient  $\rho$  in terms of expectations  $E(\cdot)$ .

$$\rho = \frac{E(ab) - E(a)E(b)}{\sqrt{(E(a^2) - E(a)^2)(E(b^2) - E(b)^2)}}$$

The only term in this formula that is influenced by the joint distribution of  $a$  and  $b$  is  $E(ab)$ . The other terms depend only on the marginals,  $a$  and  $b$ , whose descriptions are given in the problems examined in this paper. Observe from the equation that higher values of  $E(ab)$  imply higher values of  $\rho$  relative to lower values of  $E(ab)$ . This is illustrated by the fact that if  $a$  has two possible values  $a_0=9$  and  $a_1=10$ , each with a 50% chance of occurring, and  $b$  has the same distribution for its possible values  $b_0$  and  $b_1$ , then if  $a_0$  always co-occurs with  $b_0$ , and  $a_1$  with  $b_1$ , satisfying the intuitive concept of high correlation,  $E(ab)=(9*9+10*10)/2=90.5$ . This is higher than if  $a_0$  always co-occurs with  $b_1$ , and  $a_1$  with  $b_0$ , in which case  $a$  and  $b$  satisfy the intuitive

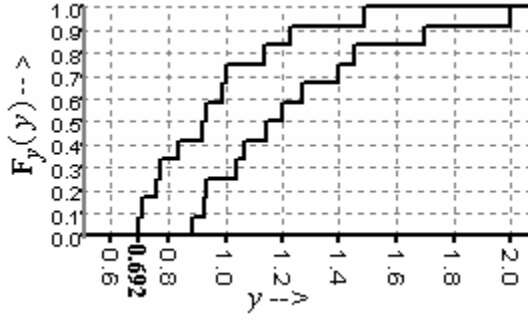
concept of low correlation and  $E(ab)=90$ . The difference is often more pronounced, as for  $a_0=-1$ ,  $a_1=1$ ,  $b_0=-1$ , and  $b_1=1$ .

If the term  $E(ab)$  is assigned a value, range, minimum, or maximum, this can be used as a constraint to augment the row and column constraints. This will tend to restrict allocation of probability masses among the interior cells of a joint distribution tableau more than the row and column constraints alone (Berleant and Zhang, forthcoming [5]). For example, suppose in Challenge Problem 3c we state that  $0.465 \leq E(ab)$ . Then the row and column constraints in Table 3 would be augmented with a new constraint, derived next.

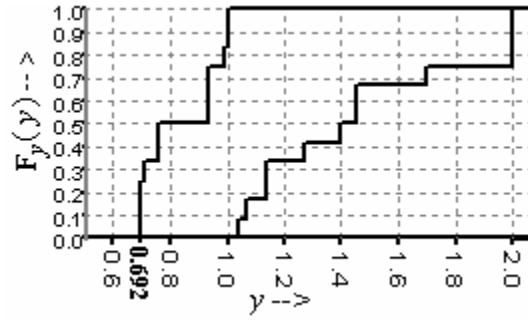
$$\begin{aligned}
0.46 &\leq \overline{\sum_{i,j} \mathbf{a}_i \cdot \mathbf{b}_j \cdot p_{ij}} && \text{(this is the constraint because we evaluate } E(ab) \text{ from the interior cells} \\
&&& \text{of a joint distribution tableau, which produces an interval that is} \\
&&& \text{consistent with the requirement that } 0.46 \leq E(ab) \text{ if any part of the} \\
&&& \text{interval is 0.465 or more)} \\
&= \sum_{i,j} \overline{\mathbf{a}_i \cdot \mathbf{b}_j \cdot p_{ij}} && \text{(the high bound of the sum is the sum of the high bounds)} \\
&= \sum_{i,j} \overline{\mathbf{a}_i} \cdot \overline{\mathbf{b}_j} \cdot p_{ij} && \text{(a number equals its high bound)} \\
&= \sum_{i,j} \overline{\mathbf{a}_i} \cdot \overline{\mathbf{b}_j} \cdot p_{ij} && \text{(the high bound of the product of non-negative intervals is the product} \\
&&& \text{of the high bounds)} \\
&= 1 \cdot 1 \cdot p_{11} + 1 \cdot 0.7 \cdot p_{21} + 1 \cdot 0.4 \cdot p_{31} + 0.7 \cdot 1 \cdot p_{12} + 0.7 \cdot 0.7 \cdot p_{22} + 0.7 \cdot 0.4 \cdot p_{32} + 0.4 \cdot 1 \cdot p_{13} \\
&+ 0.4 \cdot 0.7 \cdot p_{23} + 0.4 \cdot 0.4 \cdot p_{33} + 0.2 \cdot 1 \cdot p_{14} + 0.2 \cdot 0.7 \cdot p_{24} + 0.2 \cdot 0.4 \cdot p_{34} \geq 0.465 \\
&&& \text{(which is read directly from the joint distribution tableau in Table 3,} \\
&&& \text{top, and is the new constraint to add to Table 3, bottom).}
\end{aligned}$$

Note the assumptions  $a \geq 0$ ,  $b \geq 0$ . To remove them, see [5].

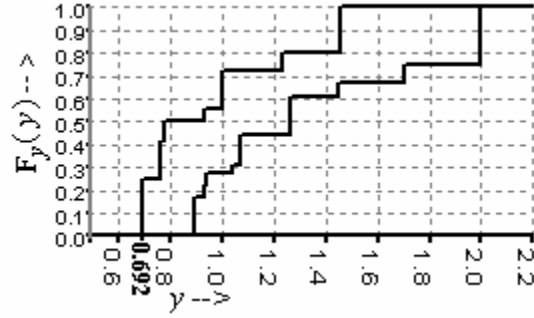
For the given marginal discretizations, the maximum value that can be attained for  $E(ab)$  for any assignment of probability masses among the interior cells of the joint distribution tableau for Problem 3c without violating the row and column constraints is 0.47. Thus the constraint  $E(ab) \geq 0.465$  enforces high correlation, which moves the envelopes closer together, excluding regions that the CDF can enter only when there is a significant chance of  $a$  being when  $b$  is low, or vice versa. This results in Figure 4 (iii). At the other extreme, the minimum value possible for  $E(ab)$  for this problem is 0.081, so the constraint  $E(ab) \leq 0.09$  enforces low correlation, excluding regions that the CDF can enter only when there is a significant chance that  $a$  and  $b$  are both low or both high. This results in Figure 4 (iv). Note the significant differences in, for example, the shape of the right tail across the conditions of independence, low and high correlation.



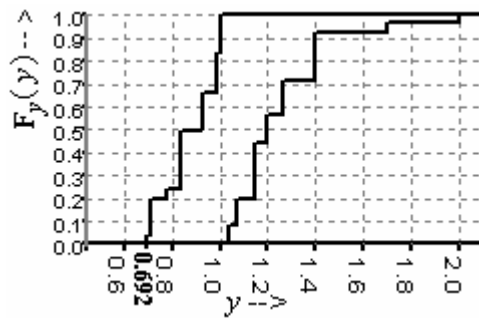
(i) Independent inputs  $a$  and  $b$ .



(ii) Unknown dependency between  $a$  and  $b$ .



(iii) High correlation between  $a$  and  $b$ .



(iv) Low correlation between  $a$  and  $b$ .

Figure 4. Solutions to Problem 3c in four variations (an additional variation will appear in Figure 14). Note that the unknown dependency condition (ii) results in envelopes that enclose those resulting from constrained dependency relationships such as those in the other three graphs.

### 2.3 Problem 4: $a$ is an interval and $b$ is a left and right envelope pair

In this problem,  $a$  is the same interval as in Problem 2. However, unlike in Problem 2, here  $b$  is a family of lognormal distributions consisting of those distributions with means in  $[0, 1]$  and standard deviations in  $[0.1, 0.5]$ . The approach taken is to convert the family into a set of intervals with associated probability masses, thereby transforming Problem 4 into one like Problem 2, and then solve as in Problem 2. Therefore this section shows:

- (1) how to convert a family of distributions into a set of intervals with associated probability masses,
- (2) how to convert a listing of intervals for  $(a+b)^a$  as in Section 2.1 into an equivalent joint distribution tableau, which is the standard format for representing problems to be solved with DEnv, and
- (3) the solution for Challenge Problem 4.

#### (1) How to convert a family of distributions into a set of intervals and probability masses.

To apply DEnv, a set of intervals and their associated probability masses must be used to represent each marginal. To convert a family of PDFs into this form, the family is first represented as a pair of left and right envelopes that enclose the CDFs of the PDFs. Then these envelopes are converted into a set of intervals and a mass for each. A method is needed that does this and has the property that the resulting intervals and associated masses will, if converted back into envelopes, produce envelopes like the ones from which the intervals and masses were derived. One such method begins by tiling the space between the envelopes with rectangles such that the left side of each rectangle is on the left envelope, and the right side of each rectangle is on the right envelope. If the envelopes have a staircase shape this may be done exactly, while in the general case this will entail discretization since the rectangles have vertical sides. The span of a rectangle over the horizontal axis defines an interval, and the bottom-to-top height of the rectangle defines the probability mass for that interval. This gives a set of intervals  $z_k$ ,  $k=1\dots K$ , and their associated masses  $p_k$ . When these are converted to envelopes according to

$$\underline{F}_z(z_0) = \sum_{k|z_k \leq z_0} p_k \text{ and } \overline{F}_z(z_0) = \sum_{k|z_k \leq z_0} p_k, \quad (3)$$

which are no more than re-subscripted forms of Equations (1), the result is envelopes identical to those from which the  $z_k$ 's were derived, or similar if the original envelopes were not staircase shaped. Figure 5 shows an example. Some issues surrounding this are discussed further in Section 3.

**(2) How to convert a 2-column table of intervals and probability masses into an equivalent joint distribution tableau.** A 2-column listing like Table 2 gives the single interval for  $a$ , lists all the intervals for  $b$ , and for each gives the result interval for  $y=(a+b)^a$ . The equivalent joint distribution tableau follows directly. Table 4 gives an example. Note how in the joint distribution tableau form, the interior cell masses  $p_{11}\dots p_{41}$  add up to 1.0 (the mass of its single marginal cell for  $a$  in the top right corner), and each mass in  $p_{11}\dots p_{41}$  equals ("adds up to") the mass of the marginal cell to its left, in accordance with the four row constraints and one column constraint implied by the tableau.

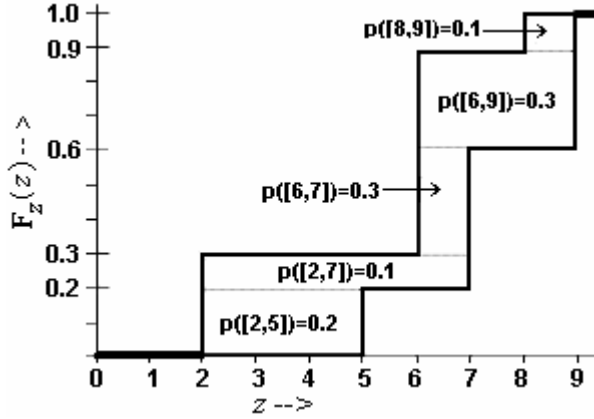


Figure 5. The left and right envelopes shown convert to the following intervals and probability masses:  $p([2,5])=0.2$ ,  $p([5,6])=0.1$ ,  $p([6,7])=0.3$ ,  $p([7,8])=0.3$ ,  $p([8,9])=0.1$ . These intervals and masses, when converted back to envelopes using Equations (3), give envelopes identical to those shown.

$b$	$y=(a+b)^a$ , given $a \in [0.1, 1]$
$[0.6, 0.8]$ $p=0.25$	$[0.96, 1.8]$ $p=0.25$
$[0.5, 0.7]$ $p=0.25$	$[0.93, 1.7]$ $p=0.25$
$[0.1, 0.4]$ $p=0.25$	$[0.76, 1.4]$ $p=0.25$
$[0, 1]$ $p=0.25$	$[0.69, 2]$ $p=0.25$

A 1-column table like Table 2 but for Problem 2c.

$b \downarrow$	$a \rightarrow$ $y \searrow$	$[0.1, 1]$ $p=1.0$
$b_1=[0.6, 0.8]$ $p=0.25$		$y_{11}=[0.96, 1.8]$ $p_{11}=0.25$
$b_2=[0.5, 0.7]$ $p=0.25$		$y_{21}=[0.93, 1.7]$ $p_{31}=0.25$
$b_3=[0.1, 0.4]$ $p=0.25$		$y_{31}=[0.76, 1.4]$ $p_{31}=0.25$
$b_4=[0, 1]$ $p=0.25$		$y_{41}=[0.69, 2]$ $p_{41}=0.25$

An equivalent joint distribution tableau.

Table 4. Example (using Problem 2c for illustration instead of Problem 4, which would require many more rows) of conversion from the simple tabular format used in Section 2.1 to the joint distribution tableau format. In fact they are almost the same.

**(3) The solution for Challenge Problem 4.** Monte Carlo simulation was used to generate envelopes enclosing the family of lognormal CDFs possible for  $b$  that were specified by Problem 4. To do this, the specified intervals for mean  $\mu$  and standard deviation  $\sigma$  were each sampled, resulting in a fully specified lognormal CDF, the height of which was then evaluated at each of a predefined set of values of  $b$ . The evaluation process was repeated for additional pairs of samples of  $\mu$  and  $\sigma$  using the same values of  $b$  and resulting in a set of CDF heights at each value of  $b$ . Then, for each value of  $b$  the highest of the CDF heights was used as a point on the left envelope and the lowest was used as a point on the right envelope. The resulting envelopes were converted to a set of 59 intervals for  $b$  as described earlier in this section. A joint distribution tableau was then constructed with 59 rows for  $b$ , and one column for the one interval provided for  $a$  by the problem statement. Then, envelopes for the cumulative distribution of  $y=(a+b)^a$  were constructed using DEnv just as in Problem 2. The results are shown in Figure 6.

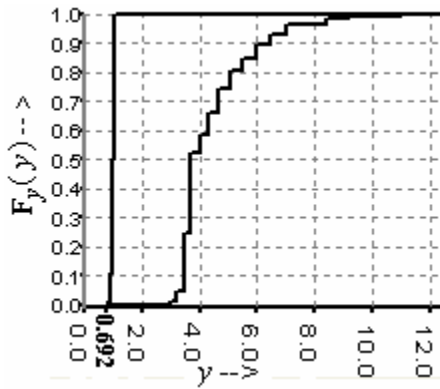


Figure 6. Solution to Problem 4.

#### 2.4 Problem 5: a set of intervals for $a$ and a set of left and right envelope pairs for $b$

The new challenge posed by this problem is dealing with not one pair of envelopes describing  $b$ , but  $n$  pairs. To handle this, we first combine the  $n$  pairs of envelopes into a single composite pair. This pair can then be converted to intervals with associated probability masses as described in Section 2.3. At that point there will be a set of intervals and their masses for  $b$ , and the set of equally weighted intervals given by the problem definition for  $a$ . The problem can then be solved like Problem 3.

**Combining envelopes.** Problems 5a-5c are similar, in that each states three left-and-right pairs of envelopes of equal credibility. Therefore we assume a total cumulation of 0.333 for  $b$  over  $-\infty$  to  $\infty$  for each of the three envelope pairs. To do this, each of the three pairs was normalized to reach a height of 0.333 instead of its original height of 1. Next, each pair of envelopes was converted into intervals and associated probability masses as in Section 2.3. At this point the sum of the masses of the intervals for each pair of envelopes is 0.333. Finally the entire set  $J$  of all of the intervals from the three pairs of envelopes, with a combined mass of 1.0, were converted to a new combined pair of envelopes in accordance with Equations (3).

Having obtained a single, combined envelope pair, an equivalent set of intervals and associated probability masses can be derived as described in Section 2.3 (or set  $J$  may simply be used), forming the marginal for  $b$ . The given intervals for  $a$  are used to form the marginal for  $a$ . At this point the problem can be solved like Problem 3 in Section 2.2.



To solve Problem 5a we first used Monte Carlo simulation as in Section 2.3 to obtain a pair of envelopes for each source of information about  $b$ . These were combined as just described. The combined pair (see Figure 7) was converted into a set of intervals and associated probability masses as in Section 2.3, which was used as a marginal for a joint distribution tableau from which the result envelopes were derived. Problems 5b and 5c were also solved this way. Figures 8-10 shows the results when  $a$  and  $b$  are independent, as well as when dependency is unknown. Note the scalloped envelopes in Figure 10, caused by gaps between the intervals given for  $a$  and for standard deviation of  $b$  [23].

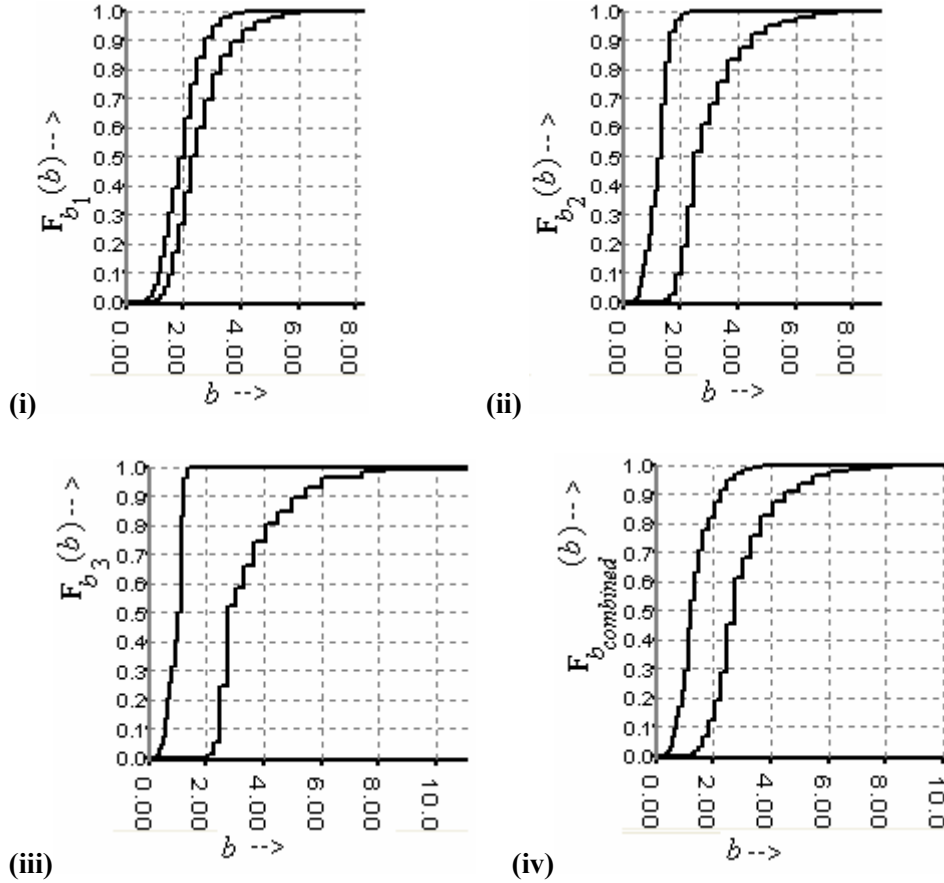


Figure 7. Three sources of information (i)-(iii) about  $b$  given for Problem 5a (see [23]) and their combination (iv). Information  $F_{b_1}$  contains PDFs with (real-valued) means in  $[0.6, 0.8]$  and (real-valued) standard deviations in  $[0.3, 0.4]$ ; (ii)  $F_{b_2}$  contains PDFs with means in  $[0.2, 0.9]$  and standard deviations in  $[0.2, 0.45]$ ; (iii)  $F_{b_3}$  contains PDFs with means in  $[0, 1]$  and standard deviations in  $[0.1, 0.5]$ . Weighting each information equally and combining yields the envelopes shown in (iv).

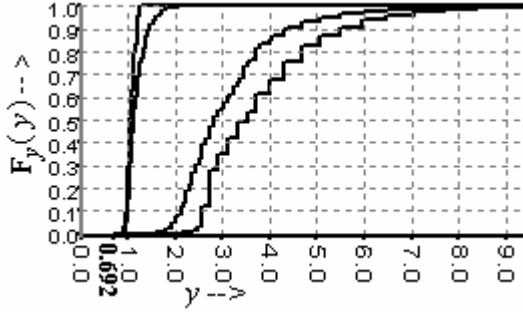


Figure 8. Results for Problem 5a when  $a$  and  $b$  are independent of each other (left and right nested envelopes) and when their dependency is unknown (left and right enclosing envelopes). Information about  $a$  is the three independent, equally weighted, nested intervals  $[0.5, 0.7]$ ,  $[0.3, 0.8]$ , and  $[0.1, 1]$ ;  $b$  is as in Figure 7 (iv).

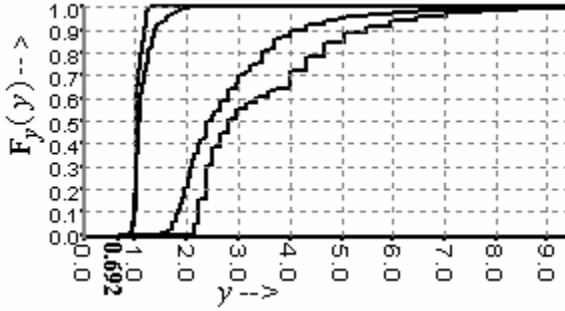


Figure 9. Results for Problem 5b when  $a$  and  $b$  are independent (nested envelopes) and when their dependency is unknown (enclosing envelopes). Information about  $a$  is three independent, equally weighted, overlapping intervals. Information about  $b$  is similar to that given in Problem 5a, except that the intervals for the means and standard deviations are different [23].

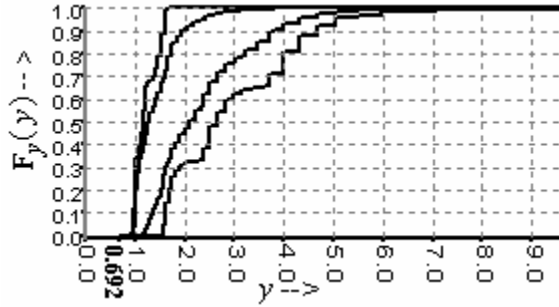
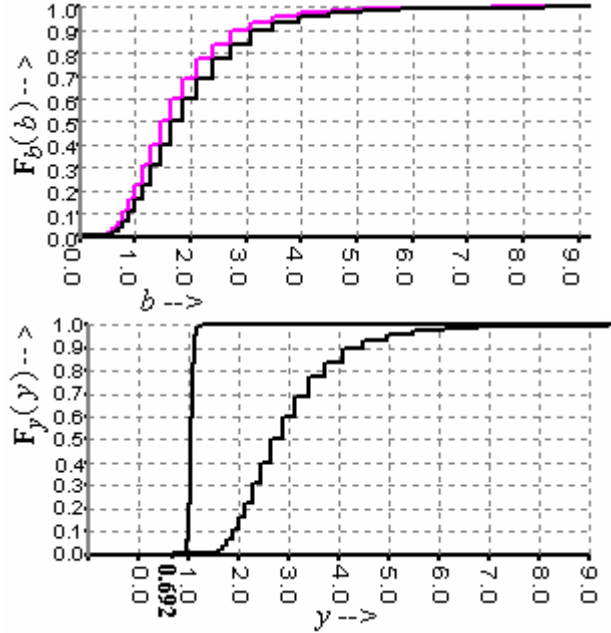


Figure 10. Results for Problem 5c when  $a$  and  $b$  are independent (nested envelopes) and when their dependency is unknown (enclosing envelopes). Information about  $a$  and  $b$  are similar in format to that given in Problems 5a and 5b.

## 2.5 Problem 6: an interval for $a$ and a distribution for $b$

When  $b$  is represented using envelopes, as in Challenge Problem 4, the envelopes may be converted into a set of intervals with associated probability masses and the solution obtained as described in Section 2.3. The salient difference in Problem 6 is that  $b$  is a single distribution rather than a family of possible distributions. To solve this, we enclosed  $b$  with left and right staircase-shaped envelopes, then convert that pair of envelopes into intervals and their probabilities, and finally solved as before. Figure 11 shows the discretization we used for  $b$ . The

southeast corners of the light left envelope touch northwest corners of the dark right envelope;  $b$  passes through those contact points. The path of  $b$  is not fully defined by the discretization over the rectangular regions between contact points, but is constrained to stay between the envelopes. Finer discretizations will constrain the path more, using more contact points and smaller rectangles between contact points. Such a discretization is safe in that it encloses rather than approximating  $b$ . Figure 11 also shows the result,  $y$ .



**Figure 11. Discretization for  $b$  in Problem 6, top, and the resulting envelopes for result  $y=(a+b)^a$ , bottom, where  $b$  is given as a lognormal PDF with mean 0.5 and standard deviation 0.5. Information about  $a$  is given as the equally credible intervals  $[0.1,0.4]$ ,  $[0.5,0.7]$ , and  $[0.8,1]$ .**

## 2.6 Problem B: the spring system

This problem involves calculating  $D_s$  from the equation

$$D_s = \frac{k}{\sqrt{(k - m\omega^2)^2 + (c\omega)^2}}.$$

Note the presence of four variables on the RHS. Each can be converted into a set of intervals and associated probabilities, the form that DEnv requires for marginals, as follows.

- (1) For  $c$ , a set of equally weighted intervals is given. This may be converted into a set of intervals and probabilities as described for Problem 2 (Section 2.1).
- (2) For  $\omega$ , a single interval-parameterized family of PDFs is given. This may be converted into a set of intervals and probabilities as described for Problem 4 (Section 2.3).
- (3) For  $m$ , a distribution is given. This may be converted to a set of intervals and probabilities as described for Problem 6 (Section 2.5).
- (4) For  $k$ , three equally credible, interval-parameterized families of PDFs are given. This set may be converted into a set of envelope pairs, these combined into one pair, and that pair converted into a set of intervals and probabilities, as described for Problem 5 (Section 2.4).

At this point, a joint distribution tableau is needed that generalizes Table 1 to 4 dimensions. It will have one marginal for each of the 4 variables. The interval for each interior cell is determined by the intervals of its four corresponding marginal cells. Thus the interval in the interior cell associated with probability mass  $p_{wxyz}$  is the range possible for  $D_s$  if  $c \in c_w$ ,  $\omega \in \omega_x$ ,  $m \in m_y$ , and  $k \in k_z$ . Interval methods can ensure that the computed interior cell intervals are not too narrow and, if too wide, are only slightly so. The four variables are given as independent, and using the standard statistical definition of independence that we have been using, the probability mass of each interior cell is the product of the masses of its four corresponding marginal cells. Once the interior cells are filled in with their intervals and probabilities, the set of interior cells may be used to generate envelopes around  $D_s$  by applying a 4-D generalization of Equations (1) or, equivalently, by numbering them consecutively and applying Equations (3). The Statool software applied to the  $(a+b)^a$  problems herein (Berleant et al. 2003 [3]) does not currently handle 4-D tableaus but an ad hoc program was written to compute the answer (Figure 12).

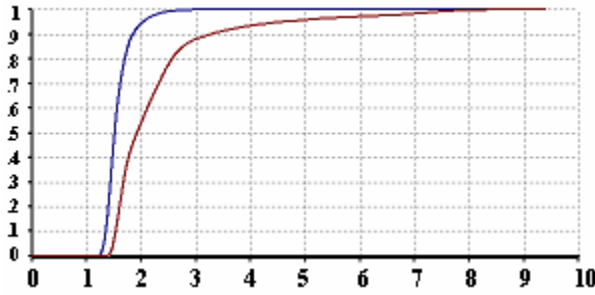


Figure 12. Envelopes around the CDF of  $D_s$  in the spring system (Challenge Problem B).

### 3 Combining Information

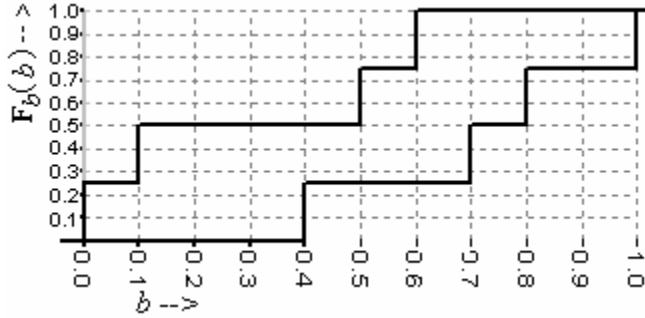
Some of the challenge problems specify  $n$  information sources of equal credibility. An important ambiguity is in the likelihood that none of the information sources are correct. This is discussed next. Then in section 3.1 we discuss information equivalence, the fact that different sets of intervals and their probabilities offered as information can be equivalent in significant ways.

#### Ambiguity in the likelihood that no source of information is correct

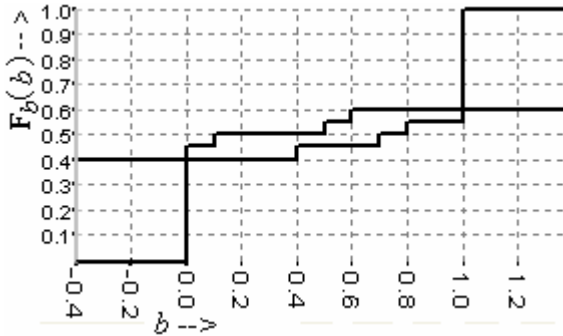
Although multiple information sources are given as having equal credibility in the challenge problems, the credibility of the true value being inconsistent with all of them (call this the *incredibility*) is unspecified. At one extreme the true value might be guaranteed to be consistent with at least one information source (zero incredibility). At the other extreme the credibilities of the information sources, though equal, might be negligible.

This ambiguity has serious implications. For example, envelopes around the CDF of  $b$  in Problem 2c can differ dramatically under different resolutions of this ambiguity. The envelopes in Figure 13 (top) were obtained under the assumption of zero incredibility, meaning the actual value of  $b$  must be binned in one of the four intervals provided by the four information sources, leading to a probability assignment of 0.25 to each. The envelopes in Figure 13 (bottom) were obtained from those in Figure 13 (top) by assigning the probability 0.05 instead of 0.25 to each bin. The four interval-valued bins then have a collective probability of  $(4) \cdot (0.05) = 0.2$ , implying an incredibility of 0.8. If we assume an incredibility of 0.8 and split it evenly into a 0.4 probability that the value of  $b$  is below all of the given intervals and a 0.4 probability that it is above all of them, the envelope pair shown will result. It consists of left and right envelopes that touch (without crossing) at two points. The middle portion, between the contact points, is like Figure 13 (top), except scaled and shifted to start at 0.4 on the vertical axis and end at 0.6. It should be noted that

solutions given earlier to challenge problems assumed at least one source of information is correct (i.e. that there is no incredibility).



Envelopes around the cumulation for  $b$  in Problem 2c, assuming no incredibility.



$p(\text{all items of information are below the actual value})=0.4$  and  
 $p(\text{all items of information are above the actual value})=0.4$

Figure 13. Envelopes around the cumulation for  $b$  in Problem 2c under two of the infinite number of possible assumptions about the probability that the actual value is not in any of the intervals given as information.

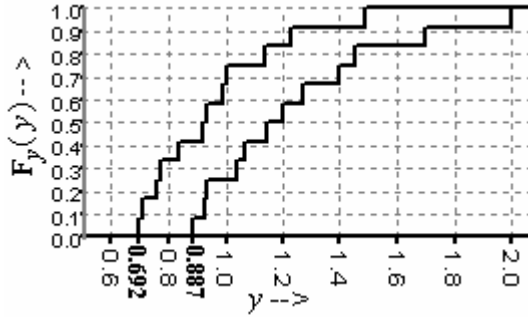
Another example of the effects of considering the possibility that no source of information is correct is that alternative solutions to Problems 3a-3c become plausible, and even arguably more plausible than the solutions given earlier. In Problems 3a-3c the three sources of information about  $a$  and four sources about  $b$  are stated to be all equally credible. This suggests that the collective credibility of all the information about  $a$  is less than the collective credibility of all the information about  $b$ , since there are fewer sources of information about  $a$ . Consequently it would make sense to model these problems with a credibility of  $\frac{1}{4}$  for each of the 3 sources of information about  $a$ , the same as for each of the 4 sources of information about  $b$ . Then a fourth possibility for  $a$ , that its value is given by some unstated or unknown information, will also have a credibility of  $\frac{1}{4}$ . This meets the requirement that all 7 information sources have equal credibility. If the  $\frac{1}{4}$  incredibility for  $a$  is spread over, for example,  $[1, 1000]$  then modeling credibility with probability implies that the CDF of  $a$  is  $\frac{3}{4}$  at  $a=1$  because of the given information sources, rising to 1 at  $a=1000$ . Consequently the CDF of  $y=(a+b)^a$  would reach 1 only at  $y=(1000+1)^{1000}$ .

Because all of the challenge problems involving  $(a+b)^a$  specified either  $[0.1, 1]$  or a subset thereof as information about  $a$ , we used that interval instead of  $[0, 1000]$  as the domain of  $a$ , and explored the implications of a  $\frac{1}{4}$  incredibility for  $a$  in the context of Problem 3c. This incredibility was divided equally between the two intervals within the  $[0.1, 1]$  domain that were not covered by any

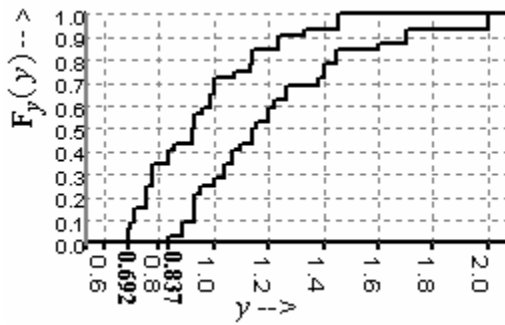
of the three intervals given by the three information sources. This means that intervals for  $a$  received the following probability assignments:

$$\begin{aligned}
 p(a \in [0.1, 0.4]) &= 0.25 \\
 p(a \in (0.4, 0.5]) &= 0.125 \\
 p(a \in [0.5, 0.7]) &= 0.25 \\
 p(a \in (0.7, 0.8]) &= 0.125 \\
 p(a \in [0.8, 1.0]) &= 0.25
 \end{aligned} \tag{3}$$

(where square brackets designate closed intervals and round brackets designate open intervals, although whether any interval is open or closed does not affect the conclusions in this example). Figure 14 (top) shows the resulting envelopes for  $y=(a+b)^a$  under the no incredibility condition, while Figure 14 (bottom) shows envelopes obtained under the incredibility scenario of Equations (3). Numerous differences exist between the two results. One is that the heights of the right envelopes differ noticeably at  $y=1.8$ . Another is that under the incredibility scenario of Equations (3) it is possible that  $a \in [0.4, 0.5]$  and  $b \in [0, 0.2]$ , in which case  $y$  can be no higher than 0.837, so the right envelope rises at 0.837. In contrast the right envelope under the no incredibility scenario rises starting at a higher number, 0.887.



No incredibility in information about  $a$ .



Result when probabilities are assigned to intervals for  $a$  per Equations (3).

Figure 14. Solutions to Problem 3c under two interpretations of information about  $a$ . In both,  $a$  and  $b$  were assumed independent.

### 3.1 Information equivalence

**Lemma 1.** It is possible for two different sets of intervals and associated probabilities to have identical envelopes.

**Proof.** By example (Figure 15).  $\square$

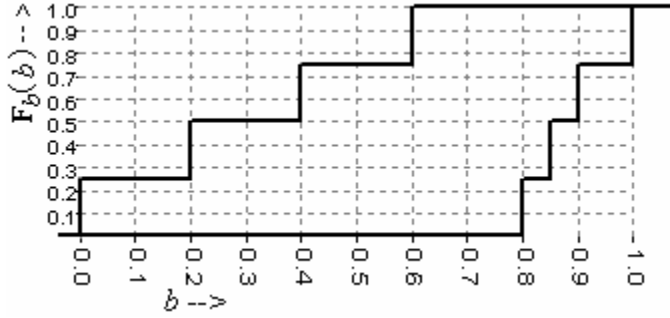


Figure 15. Envelopes around the cumulation of  $b$  in Problem 2a. They can be generated by Equations (3) from either the pieces of information given in [23], namely  $[0.6, 0.8]$ ,  $[0.4, 0.85]$ ,  $[0.2, 0.9]$ , and  $[0, 1]$ , or alternatively from the pieces of information  $[0, 0.8]$ ,  $[0.2, 0.85]$ ,  $[0.4, 0.9]$ , and  $[0.6, 1]$ .

**Lemma 2.** Consider two values  $c$  and  $d$ , each described with a different set of interval-valued pieces of information such that the envelopes computed by Equations (3) around the CDF for  $c$  are identical to those computed for  $d$ . Then, given function  $g$  and information about value  $e$ , the best-possible enclosing envelopes around the CDF of value  $g(c, e)$  are *not* necessarily identical to those around the CDF of value  $g(d, e)$ .

**Proof.** We will find a unary function  $g_1(\cdot)$  for which the envelopes around the CDF of value  $g_1(c)$  are not the same as the envelopes for  $g_1(d)$ . Since a unary function  $g_1(\cdot)$  is convertible to an equivalent binary function  $g(\cdot, \cdot)$  for which the second argument either does not affect its value or does not affect its value significantly, the lemma will be proved.

Consider function  $g_1(\cdot)$  shown in Figure 16. Let the CDF of  $c$  be partially defined by intervals  $[1, 4]$  and  $[2, 3]$ , each with probability 0.5. Let the CDF of  $d$  also be partially defined, in this case with the intervals  $[1, 3]$  and  $[2, 4]$ , each with probability 0.5. The envelopes around  $c$  are identical to those around  $d$ . This can be seen by applying Equations (3), which gives the same envelopes in both cases because both sets of information have the same low bounds and associated probabilities, and the same high bounds and associated probabilities. Although which low bound is in the same interval as which high bound differs for  $c$  and  $d$ , this does not affect the calculations of Equations (3).

Although the envelopes for  $c$  and  $d$  are identical, the envelopes for  $g_1(c)$  and  $g_1(d)$  are *not* identical. Figure 16 shows that  $g_1([1, 3]) = [6, 8] \neq [6, 7] = g_1([2, 3])$  for interval extension  $g_1(\cdot)$  of real function  $g_1(\cdot)$ , while  $g_1([1, 4]) = [6, 9] = g_1([2, 4])$ . Consequently applying Equations (3) to the information about  $g_1(c)$ , intervals  $[6, 7]$  and  $[6, 9]$  each with probability 0.5, gives a different pair of envelopes than applying Equations (3) to the information about  $g_1(d)$ , intervals  $[6, 8]$  and  $[6, 9]$  each with probability 0.5. Readers may enjoy verifying this for themselves.  $\square$

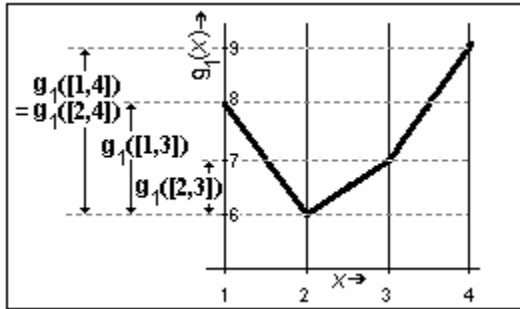


Figure 16. A function  $g_1(x)$  and projections of 4 intervals for  $x$  from the horizontal to the vertical axis. For example  $g_1(x)$  is within the range  $[6, 8]$  when  $x$  is within  $[1, 3]$  (with the 6 occurring when  $x=2$ ).

So far this section has shown that two uncertain values described by different sets of intervals and their probabilities can have identical envelopes. A function was found that takes different inputs with identical envelopes, and gives outputs with differing envelopes. On the other hand some other functions, when given different inputs with identical envelopes, will be guaranteed to output identical envelopes. Although not a full characterization of such envelope-preserving functions, the following holds.

**Theorem 1: envelope invariance for monotonic functions.** Let the available information about value  $c$  be described by a set of intervals and their associated probability masses, and similarly for  $d$ . Then Equations (3) may be applied to get envelopes around the CDFs for  $c$  and for  $d$ . If

- (i) the envelopes for  $c$  are identical to those for  $d$ ,
- (ii) the set of intervals and associated probability masses for  $c$  differs from the set for  $d$ ,
- (iii)  $g(\cdot, \cdot)$  is increasing in both arguments, and
- (iv)  $c$  and  $d$  are statistically independent in the usual sense,

then the envelopes around the CDF for value  $g(c, \cdot)$  are identical to the envelopes around the CDF for value  $g(d, \cdot)$ .

The strategy for establishing the theorem is to dissociate the effects on the result envelopes of the interval low bounds in the marginals from the effects of the high bounds. Then for constructing result envelopes it doesn't matter which marginal interval low bounds are in the same interval with which marginal interval high bounds, so that different information can lead to the same result envelopes.

- First consider point (i). Since  $c$  and  $d$  have identical envelopes, each vertical line segment in the left envelope of  $c$  has the same location on the horizontal axis and the same bottom-to-top length as a vertical line segment in the left envelope of  $d$ , and vice versa. Equations (3) imply that the horizontal axis location of each vertical line segment in the left envelope of  $c$  equals the low bound of interval(s) in the discretization of  $c$ , and the pooled mass associated with the interval(s) is the length of the segment. The same is true of the left envelope for  $d$ , and also for the right envelopes of both  $c$  and  $d$  except with reference to the interval high bounds instead of their low bounds.
- Next consider point (ii). This can be the case (by Lemma 1).
- Now consider point (iii). If  $g(x,y)$  increases monotonically in  $x$  and  $y$ , then the range of  $g(x,y)$  over the rectangle  $[x_l, x_h] \times [y_l, y_h]$  is  $[g(x_l, y_l), g(x_h, y_h)]$ . This is because  $g(x,y)$ , being monotonically increasing, has no local minima or maxima. Thus its minimum and maximum over a rectangle are at the southwest and northeast corners respectively. Each interior cell in a joint distribution tableau is determined by two marginal cells (Table 1) whose intervals define a rectangle. Therefore the interval in each interior cell has its low bound determined by the low bounds of its corresponding marginal cell intervals. Consequently the locations on the horizontal axis of the vertical line segments of the left envelope of  $g(x,y)$ , because they are determined by the low bounds of the interior cell intervals per Equations (1), are ultimately derived from the low bounds of the marginal intervals. The situation is analogous for high bounds and the right envelope. Therefore which marginal interval low bounds occur in the same interval with which high bounds is irrelevant in determining the *horizontal axis locations of vertical line segments in the result envelopes*.
- Finally, consider point (iv). In the case of independence, the probability mass of each interior cell in a joint distribution tableau is the product of the masses of its two corresponding marginal cells, and the intervals associated with these masses have no



effect. Thus the interval high bounds do not directly affect the mass computations when computing the left envelope. Similarly, interval low bounds do not directly affect the mass computations for the right envelope. It is these mass computations that determine the lengths of the vertical line segments in the envelopes. Therefore which marginal interval low bounds occur in the same interval with which high bounds is irrelevant in determining the *lengths of vertical line segments in the result envelopes*.

Thus neither the horizontal axis locations of the vertical line segments in the left envelope nor their lengths are affected by any interval high bounds in the tableau, only low bounds. Similarly, neither the horizontal axis locations of the vertical line segments in the right envelope nor their lengths are affected by any interval low bounds in the tableau, only high bounds. Therefore which marginal interval high bounds occur in the same intervals with which marginal interval low bounds does not affect the envelopes and the theorem is established.  $\square$

*Further comments.* The irrelevance of which high bound is associated with which low bound explains why the envelopes of Figure 15 are implied by both of the two different information sets noted in that figure.

Further work is needed to better understand equivalence of information sets in the presence of unknown dependency and correlation constraints, as well as in cases where  $g(x,y)$  is not monotonically increasing in  $x$  and  $y$ . Some recent work relevant to such questions includes Hall and Lawry (this issue [15]), Ferson and Kreinovich [12], and Kreinovich et al. 2001 [18].

## 4 Conclusion

The challenge problems of Oberkampff et al. [23] provide a valuable opportunity to compare techniques and their implementations for computing functions of random variables whose samples are expressed using intervals, distributions, or families of distributions. Monte Carlo techniques form a conceptually graspable class of algorithms but give results complicated by random noise. Probabilistic Arithmetic, random set theory, and joint distribution tableaus, which can be manipulated by the DEnv algorithm, are the most widely reported alternatives. Of these three, joint distribution tableaus do not require an understanding of copulas or random sets, so are arguably easier to understand. DEnv can compute functions of random variables when they are either: (i) independent of each other, (ii) have some other specific dependency relationship, (iii) have an unknown dependency relationship, or (iv) have a dependency that is partially characterized. Results provided by DEnv are consistent with those provided by the other techniques to date. Additional work is needed to achieve further advances in such directions as more flexible handling of partially characterized dependency. Ultimately, crossing the bridge from 2<sup>nd</sup> order probabilistic results to decisions is likely to be a key factor in their achieving widespread use.

## Acknowledgements

The authors are grateful to Lizhi Xie for building an initial version of Statool suitable for continued development, to Gerald Sheblé for motivating advances in Statool through extensive discussions on the needs of applications in the electric power industry, especially power systems economics, and to Scott Ferson and Vladik Kreinovich for valuable discussions and well-timed encouragement. We also acknowledge support from PSERC to G. Sheblé, D. Berleant, and R. Thomas.

## References

- [1] Berleant D. Automatically verified reasoning with both intervals and probability density functions. *Interval Computations* (1993 No. 2), pp. 48-70.
- [2] Berleant D and Goodman-Strauss C. Bounding the results of arithmetic operations on random variables of unknown dependency using intervals. *Reliable Computing* 4 (2) (1998), pp. 147-165.
- [3] Berleant D, Xie L, and Zhang J. Statool: a tool for Distribution Envelope Determination (DEnv), an interval-based algorithm for arithmetic on random variables. *Reliable Computing* 9 (2) (2003), pp. 91-108.
- [4] Berleant D, Zhang J, Hu R, and Sheblé G. Economic dispatch: applying the interval-based distribution envelope algorithm to an electric power problem. *SIAM Workshop on Validated Computing 2002 (Extended Abstracts)*, Toronto, May 23-25, 2002, pp. 26-31.
- [5] Berleant D and Zhang J. Using correlation to improve envelopes around derived distributions. *Reliable Computing*, accepted. [www.public.iastate.edu/~berleant](http://www.public.iastate.edu/~berleant).
- [6] Colombo AG and Jaarsma RJ. A powerful numerical method to combine random variables. *IEEE Transactions on Reliability* R-29 (2) (1980), pp. 126-129.
- [7] Couso I, Moral S, and Walley P. Examples of independence for imprecise probabilities, *1<sup>st</sup> Int. Symp. On Imprecise Probabilities and their Applications*, Ghent, Belgium, 1999. [decsai.ugr.es/~smc/isipta99/proc/proceedings.html](http://decsai.ugr.es/~smc/isipta99/proc/proceedings.html).
- [8] Ferson S. *RAMAS Risk Calc 4.0: risk assessment with uncertain numbers*. Lewis Press, Boca Raton, 2002.
- [9] Ferson S. What Monte Carlo methods cannot do. *Journal of Human and Ecological Risk Assessment* 2 (4) (1996), pp. 990-1007.
- [10] Ferson S, Ginzburg L, and Akçakaya R. Whereof one cannot speak: when input distributions are unknown, *Risk Analysis*, to appear.
- [11] Ferson S and Hajagos J. Rigorous and (often) best-possible answers to arithmetic problems. *Reliability Engineering and System Safety*, this issue.
- [12] Ferson S and Kreinovich V. Representation, elicitation, and aggregation of uncertainty in risk analysis – from traditional probabilistic techniques to more general, more realistic approaches: a survey. Manuscript, [vladik@cs.utep.edu](mailto:vladik@cs.utep.edu).
- [13] Fetz T and Oberguggenberger M. Propagation of uncertainty through multivariate functions in the framework of sets of probability measures. *Reliability Engineering and System Safety*, this issue.
- [14] Frank MJ, Nelsen RB, and Schweizer B. Best-possible bounds for the distribution of a sum – a problem of Kolmogorov. *Probability Theory and Related Fields* 74 (1987), pp. 199-211.
- [15] Hall J and J Lawry. Generation, combination and extension of random set approximations to coherent lower and upper probabilities. *Reliability Engineering and System Safety*, this issue.
- [16] Ingram GE, Welker EL, and Herrmann CR. Designing for reliability based on probabilistic modeling using remote access computer systems. *Proc. 7th Reliability and Maintainability Conference*, American Society of Mechanical Engineers, 1968, pp. 492-500.
- [17] Kaplan S. On the method of discrete probability distributions in risk and reliability calculations, applications to seismic risk assessment. *Risk Analysis* 1 (3) (1981), pp. 189-196.
- [18] Kreinovich V, Langrand C, and Nguyen HT. Combining fuzzy and probabilistic knowledge using belief functions. *Proc. 2<sup>nd</sup> Vietnam-Japan Bilateral Symposium on Fuzzy Systems and Applications 2001*, Hanoi, Dec. 7-8, pp. 191-198.

- [19] Lodwick W and Jamison KD. Estimating and validating the cumulative distribution of a function of random variables: toward the development of distribution arithmetic. *Reliable Computing* 9 (2) (2003), pp. 127-141.
- [20] Moore RE. Risk analysis without Monte Carlo methods. *Freiburger Intervall-Berichte*, 84/1, 1984, pp. 1-48.
- [21] Nelsen RB. *An Introduction to Copulas*. Lecture Notes in Statistics, Vol. 139, 1999, Springer-Verlag.
- [22] Neumaier A. Clouds, fuzzy sets and probability intervals. Submitted.  
[www.mat.univie.ac.at/~neum/ms/papers.html](http://www.mat.univie.ac.at/~neum/ms/papers.html).
- [23] Oberkampf WL, Helton JC, Joslyn CA, Wojtkiewics SF, and Ferson S. Challenge problems: uncertainty in system response given uncertain parameters. *Reliability Engineering and System Safety*, this issue.
- [24] Red-Horse J. and Benjamin AS. A probabilistic approach to uncertainty quantification with limited information. *Reliability Engineering and System Safety*, this issue.
- [25] Regan H, Ferson S, and Berleant D. Equivalence of five methods for bounding uncertainty. *Int. Journal of Approximate Reasoning*, accepted pending revisions. Draft at [www.public.iastate.edu/~berleant](http://www.public.iastate.edu/~berleant).
- [26] Sandia National Laboratory. *Epistemic Uncertainty Workshop*, August 6-7, 2002, Albuquerque. Presentations and papers are at [www.sandia.gov/epistemic/](http://www.sandia.gov/epistemic/).
- [27] Sheblé G and Berleant D. Bounding the composite value at risk for energy service company operation with DEnv, an interval-based algorithm. *SIAM Workshop on Validated Computing 2002 (Extended Abstracts)*, Toronto, May 23-25, 2002, pp. 166-171.
- [28] Springer MD. *The Algebra of Random Variables*, John Wiley and Sons, New York, 1979.
- [29] Tonon F. Using random set theory to propagate epistemic uncertainty through a mechanical system. *Reliability Engineering and System Safety*, this issue.
- [30] Williamson RC and Downs T. Probabilistic Arithmetic I: numerical methods for calculating convolutions and dependency bounds. *International Journal of Approximate Reasoning* 4 (1990), pp. 89-158.