

Adding Unimodality or Independence Makes Interval Probability Problems NP-Hard

Daniel J. Berleant¹, Olga Kosheleva²,
Vladik Kreinovich², and Hung T. Nguyen³

¹Department of Electrical and Computer Engineering
Iowa State University, Ames, IA 50011, USA
berleant@iastate.edu

²NASA Pan-American Center for
Earth and Environmental Studies (PACES)
University of Texas, El Paso, TX 79968, USA
olgak@utep.edu, vladik@utep.edu

³Department of Mathematical Sciences
New Mexico State University, Las Cruces, NM 88003, USA
hunguyen@nmsu.edu

Abstract

In many real-life situations, we only have partial information about probabilities. This information is usually described by bounds on moments, on probabilities of certain events, etc. – i.e., by characteristics $c(p)$ which are linear in terms of the unknown probabilities p_j . If we know interval bounds on some such characteristics $\underline{a}_i \leq c_i(p) \leq \bar{a}_i$, and we are interested in a characteristic $c(p)$, then we can find the bounds on $c(p)$ by solving a linear programming problem.

In some situations, we also have additional conditions on the probability distribution – e.g., we may know that the two variables x_1 and x_2 are independent, or that for each value of x_2 , the corresponding conditional distribution for x_1 is unimodal. We show that adding each of these conditions makes the corresponding interval probability problem NP-hard.

1 Introduction

Interval probability problems can be often reduced to linear programming (LP). In many real-life situations, in addition to the *intervals* $[\underline{x}_i, \bar{x}_i]$

of possible values of the unknowns x_1, \dots, x_n , we also have partial information about the *probabilities* of different values within these intervals.

This information is usually given in terms of bounds on the standard characteristics $c(p)$ of the corresponding probability distribution p , such as the k -th moment $M_k \stackrel{\text{def}}{=} \int x^k \cdot \rho(x) dx$ (where $\rho(x)$ is the probability density), the values of the cumulative distribution function (cdf) $F(t) \stackrel{\text{def}}{=} \text{Prob}(x \leq t) = \int_{-\infty}^t \rho(x) dx$ of some of the variables, etc. Most of these characteristics are linear in terms of $\rho(x)$ – and many other characteristics like central moments are combinations of linear characteristics: e.g., variance V can be expressed as $V = M_2 - M_1^2$.

A typical practical problem is when we know the ranges of some of these characteristics $\underline{a}_i \leq c_i(p) \leq \bar{a}_i$, and we want to find the range of possible values of some other characteristic $c(p)$. For example, we know the bounds on the marginal cdfs for the variables x_1 and x_2 , and we want to find the range of values of the cdf for $x_1 + x_2$.

In such problems, the range of possible values of $c(p)$ is an interval $[\underline{a}, \bar{a}]$. To find \underline{a} (correspondingly, \bar{a}), we must minimize (correspondingly, maximize) the linear objective function $c(p)$ under linear constraints — i.e., solve a linear programming (LP) problem; see, e.g., [16, 17, 18, 21].

Other simple examples of linear conditions include bounds on the values of the density function $\rho(x)$; see, e.g., [15].

Comment. Several more complex problems can also be described in LP terms. For example, when we select a new strategy for a company (e.g., for an electric company), one of the reasonable criteria is that the expected monetary gain should be not smaller than the expected gain for a previously known strategy. In many cases, for each strategy, we can estimate the probability of different production values – e.g., the probability $F(t) = \text{Prob}(x \leq t)$ that we will produce the amount $\leq t$. However, the utility $u(t)$ corresponding to producing t depends on the future prices and is not well known; therefore, we cannot predict the exact value of the expected utility $\int u(x) \cdot \rho(x) dx$. One way to handle this situation is require that for *every* monotonic utility function $u(t)$, the expected utility under the new strategy – with probability density function (pdf) $\rho(x)$ and cdf $F(x)$ – is larger than or equal to the expected utility under the old strategy – with pdf $\rho_0(x)$ and cdf $F_0(x)$: $\int u(x) \cdot \rho(x) dx \geq \int u(x) \cdot \rho_0(x) dx$. This condition is called *first order stochastic dominance*. It is known that this condition is equivalent to the condition that $F(x) \leq F_0(x)$ for all x .

Indeed, the condition is equivalent to

$$\int_0^t u(x) \cdot (\rho(x) - \rho_0(x)) dx \geq 0.$$

Integrating by part, we conclude that

$$- \int_0^t u'(x) \cdot (F(x) - F_0(x)) dx \geq 0;$$

since $u(x)$ is non-decreasing, the derivative $u'(x)$ can be an arbitrary non-negative function; so, the above condition is indeed equivalent to $F(x) \leq F_0(x)$ for all x .

Each of these inequalities is linear in terms of $\rho(x)$ – so, optimizing a linear objective function under the constraints $F(x) \geq F_0(x)$ is also a LP problem.

This requirement may be too restrictive; in practice, preferences have the property of risk aversion: it is better to gain a value x with probability 1 than to have either 0 or $2x$ with probability $1/2$. In mathematical terms, this condition means that the corresponding utility function $u(x)$ is concave. It is therefore reasonable to require that for all such *risk aversion* utility functions $u(x)$, the expected utility under the new strategy is larger than or equal to the expected utility under the old strategy. This condition is called *second order stochastic dominance* (see, e.g., [8, 9, 19, 20]), and it known to be equivalent to the condition that $\int_0^t F(x) dx \leq \int_0^t F_0(x) dx$.

Indeed, the condition is equivalent to

$$\int_0^t u(x) \cdot (\rho(x) - \rho_0(x)) dx \geq 0$$

for every concave function $u(x)$. Integrating by part twice, we conclude that

$$\int_0^t u''(x) \cdot \left(\int_0^x F(z) dz - \int_0^x F_0(z) dz \right) dx \geq 0.$$

Since $u(x)$ is concave, the second derivative $u''(x)$ can be an arbitrary non-positive function; so, the above condition is indeed equivalent to $\int_0^t F(x) dx \leq \int_0^t F_0(x) dx$ for all t .

The cdf $F(x)$ is a linear combination of the values $\rho(x)$; thus, its integral $\int F(x) dx$ is also linear linear in $\rho(x)$, and hence the above condition is still linear in terms of the values $\rho(x)$. Thus, we again have a LP problem; for details, see [2].

Most of the corresponding LP problems can be efficiently solved.

Theoretically, some of these LP problems have infinitely many variables $\rho(x)$, but in practice, we can discretize each coordinate and thus, get a LP problem with finitely many variables.

There are known efficient algorithms and software for solving LP problems with finitely many variables. These algorithms require polynomial time ($\leq n^k$) to solve problems with $\leq n$ unknowns and $\leq n$ constraints; these algorithms are actively used in imprecise probabilities; see, e.g., [1, 1, 4, 5, 6, 7].

For example, for the case of two variables x_1 and x_2 , we may know the probabilities $p_i = p(x_1 \in [i, i + 1])$ and $q_j = p(x_2 \in [j, j + 1])$ for finitely many intervals $[i, i + 1]$. Then, to find the range of possible values of, e.g.,

$$\text{Prob}(x_1 + x_2 \leq k),$$

we can consider the following linear programming problem: the unknowns are

$$p_{i,j} \stackrel{\text{def}}{=} p(x_1 \in [i, i+1] \& x_2 \in [j, j+1]),$$

the constraints are $p_{i,j} \geq 0$, $p_{i,1} + p_{i,2} + \dots = p_i$, $p_{1,j} + p_{2,j} + \dots = q_j$, and the objective function is $\sum_{i,j:i+j \leq k} p_{i,j}$.

Comment. The only LP problems for which there may not be an efficient solution are problems involving a large amount of variables v . If we discretize each variable into n intervals, then overall, we need n^v unknowns p_{i_1, i_2, \dots, i_v} ($1 \leq i_1 \leq n$, $1 \leq i_2 \leq n$, \dots , $1 \leq i_v \leq n$) to describe all possible probability distributions. When v grows, the number of unknowns grows exponentially with v and thus, for large v , becomes unrealistically large.

It is known (see, e.g., [13]) that this exponential increase in complexity is inherent to the problem: e.g., for v random variables x_1, \dots, x_v with known marginal distributions, the problem of finding the exact bounds on the cdf for the sum $x_1 + \dots + x_v$ is NP-hard.

Beyond LP. There are important practical problems which lie outside LP. One example is problems involving *independence*, when constraints are linear in $p(x, y) = p(x) \cdot p(y)$ and thus, bilinear in $p(x)$ and $p(y)$. In this paper, we show that the corresponding range estimation problem is NP-hard.

Another example of a condition which cannot be directly described in terms of LP is the condition of *unimodality*. For a one-variable distribution with probabilities p_1, \dots, p_n , unimodality means that there exists a value m (“mode”) such that p_i increase (non-strictly) until m and then decreases after m :

$$p_1 \leq p_2 \leq \dots \leq p_{m-1} \leq p_m \geq p_{m+1} \geq \dots \geq p_{n-1} \geq p_n.$$

When the location of the mode is known, we get several linear inequalities, so we can still use efficient techniques such as LP; see, e.g., [10, 22].

For a 1-D case, if we do not know the location of the mode, we can try all n possible locations and solve n corresponding LP problems. Since each LP problem requires a polynomial time to run, running n such problems still requires a polynomial time.

In the 2-D case, it is reasonable to consider the situation when, e.g., for every value of x_2 , the corresponding conditional distribution for x_1 is unimodal. In this case, to describe this as a LP problem, we must select a mode for every x_2 . If there are n values of x_2 , and at least 2 possible choices of mode location, then we get an exponential amount of 2^n possible choices. In this paper, we show that this problem is also NP-hard – and therefore, that, unless P=NP, no algorithm can solve it in polynomial time.

Comment.

- Some of the results presented in this paper have appeared in our conference paper [3].
- Other possible restrictions on probability may involve bounds on the *entropy* of the corresponding probability distributions; such problems are also, in general, NP-hard [12].

2 Adding Unimodality Makes Interval Probability Problems NP-Hard

Definition 1 Let $n_1 > 0$ and $n_2 > 0$ be given integers.

- By a probability distribution, we mean a collection of real numbers $p_{i_1, i_2} \geq 0$, $1 \leq i_1 \leq n_1$, and $1 \leq i_2 \leq n_2$, such that $\sum_{i_1=1}^{n_1} \sum_{i_2=1}^{n_2} p_{i_1, i_2} = 1$.

- We say that the distribution p_{i_1, i_2} is unimodal in the 1st variable (or 1-unimodal, for short) if for every i_2 from 1 to n_2 , there exists a value m such that p_{i_1, i_2} grows with i_1 for $i_1 \leq m$ and decreases with i_1 for $i_1 \geq m$:

$$p_{1, i_2} \leq p_{2, i_2} \leq \dots \leq p_{m, i_2} \geq p_{m+1, i_2} \geq \dots \geq p_{n_1, i_2}.$$

- We say that the distribution p_{i_1, i_2} is unimodal in the 2nd variable (or 2-unimodal, for short) if for every i_1 from 1 to n_1 , there exists a value m such that p_{i_1, i_2} grows with i_2 for $i_2 \leq m$ and decreases with i_2 for $i_2 \geq m$:

$$p_{i_1, 1} \leq p_{i_1, 2} \leq \dots \leq p_{i_1, m} \geq p_{i_1, m+1} \geq \dots \geq p_{i_1, n_2}.$$

- We say that the distribution p_{i_1, i_2} is unimodal if it is both 1-unimodal and 2-unimodal.
- By a linear constraint on the probability distribution, we mean the constraint of the type $\underline{b} \leq \sum_{i_1=1}^{n_1} \sum_{i_2=1}^{n_2} b_{i_1, i_2} \cdot p_{i_1, i_2} \leq \bar{b}$ for some given values \underline{b} , \bar{b} , and b_{i_1, i_2} .
- By an interval probability problem under 1-unimodality constraint, we mean the following problem: given a find list of linear constraints, check whether there exists a 1-unimodal distribution which satisfies all these constraints.
- By an interval probability problem under unimodality constraint, we mean the following problem: given a find list of linear constraints, check whether there exists a unimodal distribution which satisfies all these constraints.

Theorem 1 Interval probability problem under 1-unimodality constraint is NP-hard.

Comment. This clearly means that the interval probability problem under 2-unimodality constraint is also NP-hard.

Theorem 2 *Interval probability problem under unimodality constraint is NP-hard.*

Comment. So, under the unimodality constraints, even checking whether a system of linear constraints is consistent – i.e., whether the range of a given characteristic is empty – is computationally difficult (NP-hard).

Proof of Theorem 1. We will show that if we can check, for every system of linear constraints, whether this system is consistent or not under unimodality, then we would be able to solve a *partition* problem which is known to be NP-hard [11, 14]. The partition problem consists of the following: given n positive integers s_1, \dots, s_n , check whether exist n integers $\varepsilon_i \in \{-1, 1\}$ for which $\varepsilon_1 \cdot s_1 + \dots + \varepsilon_n \cdot s_n = 0$.

Indeed, for every instance of the partition problem, we form the following system of constraints: $n_1 = 3$, $n_2 = n$,

- $p_{2,i_2} = 0$ for every $i_2 = 1, \dots, n_2$,
- $p_{1,i_2} + p_{2,i_2} + p_{3,i_2} = 1/n$ for every $i_2 = 1, \dots, n_2$;
- $\sum_{i_2=1}^{n_2} (-s_{i_2} \cdot p_{1,i_2} + s_{i_2} \cdot p_{3,i_2}) = 0$.

Let us prove that this system is consistent if and only if the original instance of the partition problem has a solution.

“If” part. If the original instance has a solution $\varepsilon_i \in \{-1, 1\}$, then, for every i_2 from 1 to n_2 , we can take $p_{2+\varepsilon_{i_2},i_2} = 1/n$ and $p_{i_1,i_2} = 0$ for $i_1 \neq 2 + \varepsilon_{i_2}$. In other words:

- if $\varepsilon_{i_2} = -1$, then we take $p_{1,i_2} = 1/n$ and $p_{2,i_2} = p_{3,i_2} = 0$;
- if $\varepsilon_{i_2} = 1$, then we take $p_{1,i_2} = p_{2,i_2} = 0$ and $p_{3,i_2} = 1/n$.

The resulting distribution is unimodal: indeed, for each i_2 , its mode is the value $1 + \varepsilon_{i_2}$. Let us check that it satisfies all the desired constraints. It is easy to check that for every i_2 , we have $p_{2,i_2} = 0$ and $p_{1,i_2} + p_{2,i_2} + p_{3,i_2} = 1/n$. Finally, due to our choice of p_{i_1,i_2} , we conclude that $-s_{i_2} \cdot p_{1,i_2} + s_{i_2} \cdot p_{3,i_2} = \frac{1}{n} \cdot \varepsilon_{i_2} \cdot s_{i_2}$ and thus,

$$\sum_{i_2=1}^{n_2} (-s_{i_2} \cdot p_{1,i_2} + s_{i_2} \cdot p_{3,i_2}) = \frac{1}{n} \cdot \sum_{i_2=1}^{n_2} \varepsilon_{i_2} \cdot s_{i_2} = 0.$$

“Only if” part. Vice versa, let us assume that we have a unimodal distribution p_{i_1, i_2} for which all the desired constraints are satisfied. Since the distribution is unimodal, for every i_2 , there exists a mode $m_{i_2} \in \{1, 2, 3\}$ for which the values p_{i_1, i_2} increase for $i_1 \leq m_{i_2}$ and decrease for $i_1 \geq m_{i_2}$. This mode cannot be equal to 2, because otherwise, the value $p_{2, i_2} = 0$ will be the largest of the three values p_{1, i_2} , p_{2, i_2} , and p_{3, i_2} hence all three values will be 0 – which contradicts to the constraint $p_{1, i_2} + p_{2, i_2} + p_{3, i_2} = 1/n$. Thus, this mode is either 1 or 3:

- if the mode is 1, then due to monotonicity, we have $0 = p_{2, i_2} \geq p_{3, i_2}$ hence $p_{3, i_2} = p_{2, i_2} = 0$;
- if the mode is 3, then due to monotonicity, we have $p_{1, i_2} \leq p_{2, i_2} = 0$ hence $p_{1, i_2} = p_{2, i_2} = 0$.

In both case, for each i_2 , only one value of p_{i_1, i_2} is different from 0 – the value $p_{m_{i_2}, i_2}$. Since the sum of these three values is $1/n$, this non-zero value must be equal to $1/n$. If we denote $\varepsilon_i \stackrel{\text{def}}{=} m_i - 2$, then we conclude that $\varepsilon_i \in \{-1, 1\}$. For each i_2 , we have

$$-s_{i_2} \cdot p_{1, i_2} + s_{i_2} \cdot p_{3, i_2} = \varepsilon_{i_2} \cdot s_{i_2} \cdot (1/n),$$

hence from the constraint

$$\sum_{i_2=1}^{n_2} (-s_{i_2} \cdot p_{1, i_2} + s_{i_2} \cdot p_{3, i_2}) = \frac{1}{n} \cdot \sum_{i_2=1}^{n_2} \varepsilon_{i_2} \cdot s_{i_2} = 0,$$

we conclude that $\sum \varepsilon_i \cdot s_i = 0$, i.e., that the original instance of the partition problem has a solution.

The theorem is proven.

Comment. The above constraints are not just mathematical tricks, they have a natural interpretation if for x_1 , we take the values $-1, 0$, and 1 as corresponding to $i_1 = 1, 2, 3$, and for x_2 , we take the values s_1, \dots, s_n . Then:

- the constraint $p_{2, i_2} = 0$ means that $\text{Prob}(x_1 = 0) = 0$;
- the constraint $p_{1, i_2} + p_{2, i_2} + p_{3, i_2} = 1/n$ means that $\text{Prob}(x_2 = s_i) = 1/n$ for all n values s_i , and
- the constraint $\sum_{i_2=1}^{n_2} (-s_{i_2} \cdot p_{1, i_2} + s_{i_2} \cdot p_{3, i_2}) = 0$ means that the expected value of the product is 0: $E[x_1 \cdot x_2] = 0$.

So, the difficult-to-solve problem is to check whether it is possible that $E[x_1 \cdot x_2] = 0$ and $\text{Prob}(x_1 = 0) = 0$ for some unimodal distribution for which the marginal distribution on x_2 is “uniform”.

Proof of Theorem 2. To proof this result, we will reduce, to this problem, the same partition problem as in the proof of Theorem 1.

For every instance of the partition problem, we form the following system of constraints: $n_1 = 3 \cdot n$, $n_2 = n$,

- $p_{i_1, i_2} = 0$ for every $i_2 = 1, \dots, n_2$ and for every $i_2 \neq 3 \cdot i_2$ and $i_2 \neq 3 \cdot i_2 - 2$;
- $\sum_{i_1=1}^{n_1} p_{i_1, i_2} = 1/n$ for every $i_2 = 1, \dots, n_2$;
- $\sum_{i_2=1}^{n_2} (-s_{i_2} \cdot p_{3 \cdot i_2 - 2, i_2} + s_{i_2} \cdot p_{3 \cdot i_2, i_2}) = 0$.

Let us prove that this system is consistent if and only if the original instance of the partition problem has a solution.

“If” part. If the original instance has a solution $\varepsilon_i \in \{-1, 1\}$, then, for every i_2 from 1 to n_2 , we can take $p_{3 \cdot i_2 - 1 + \varepsilon_{i_2}, i_2} = 1/n$ and $p_{i_1, i_2} = 0$ for $i_1 \neq 3 \cdot i_2 - 1 + \varepsilon_{i_2}$. In other words:

- if $\varepsilon_{i_2} = -1$, then we take $p_{3 \cdot i_2 - 2, i_2} = 1/n$ and $p_{3 \cdot i_2, i_2} = 0$;
- if $\varepsilon_{i_2} = 1$, then we take $p_{3 \cdot i_2 - 2, i_2} = 0$ and $p_{3 \cdot i_2, i_2} = 1/n$.

The resulting distribution is 1-unimodal: indeed, for each i_2 , its mode is the value $3 \cdot i_2 - 1 + \varepsilon_{i_2}$. Similarly, it is 2-unimodal, because for each i_1 , only one probability p_{i_1, i_2} may be different from 0. Similarly to the proof of Theorem 1, we can check that this distribution satisfies all the desired constraints.

“Only if” part. Vice versa, let us assume that we have a unimodal distribution p_{i_1, i_2} for which all the desired constraints are satisfied. Since the distribution is unimodal, it is 1-unimodal, so for every i_2 , there exists a mode $m_{i_2} \in \{3 \cdot i_2 - 2, 3 \cdot i_2\}$ for which the values p_{i_1, i_2} increase for $i_1 \leq m_{i_2}$ and decrease for $i_1 \geq m_{i_2}$. Similarly to the proof of Theorem 1, we conclude that for each i_2 , only one value of p_{i_1, i_2} is different from 0 – the value $p_{m_{i_2}, i_2}$, and this non-zero value is equal to $1/n$. If we denote $\varepsilon_i \stackrel{\text{def}}{=} m_i - (3 \cdot i - 1)$, then we conclude that $\varepsilon_i \in \{-1, 1\}$. For each i_2 , we have

$$-s_{i_2} \cdot p_{3 \cdot i_2 - 2, i_2} + s_{i_2} \cdot p_{3 \cdot i_2, i_2} = \varepsilon_{i_2} \cdot s_{i_2} \cdot (1/n),$$

hence from the constraint

$$\sum_{i_2=1}^{n_2} (-s_{i_2} \cdot p_{3 \cdot i_2 - 2, i_2} + s_{i_2} \cdot p_{3 \cdot i_2, i_2}) = \frac{1}{n} \cdot \sum_{i_2=1}^{n_2} \varepsilon_{i_2} \cdot s_{i_2} = 0,$$

we conclude that $\sum \varepsilon_i \cdot s_i = 0$, i.e., that the original instance of the partition problem has a solution.

The theorem is proven.

Comment. We can get a natural interpretation of the above constraints if for $s \stackrel{\text{def}}{=} 3 \cdot \max_i s_i$, we take the values $x_2 = s \cdot i_2$ corresponding to $i_2 = 1, \dots, n_2$, and for i_1 :

- we take the value $x_1 = s \cdot i_2$ corresponding to $i_1 = 3 \cdot i_2 - 1$;
- we take the value $x_1 = s \cdot i_2 - s_{i_2}$ corresponding to $i_1 = 3 \cdot i_2 - 2$; and
- we take the value $x_1 = s \cdot i_2 + s_{i_2}$ corresponding to $i_1 = 3 \cdot i_2$.

In this interpretation, the above constraint $\sum_{i_2=1}^{n_2} (-s_{i_2} \cdot p_{3 \cdot i_2 - 2, i_2} + s_{i_2} \cdot p_{3 \cdot i_2, i_2}) = 0$ simply means that $E[x_1] = E[x_2]$.

3 Adding Independence Makes Interval Probability Problems NP-Hard

In general, in statistics, independence makes problems easier. We will show, however, that for interval probability problems, the situation is sometimes opposite: the addition of independence assumption turns easy-to-solve problems into NP-hard ones.

Definition 2 *Let $n_1 > 0$ and $n_2 > 0$ be given integers.*

- *By an independent probability distribution, we mean a collection of real numbers $p_i \geq 0$, $1 \leq i \leq n_1$, and q_j , $1 \leq j \leq n_2$, such that $\sum_{i=1}^{n_1} p_i = \sum_{j=1}^{n_2} q_j = 1$.*

- *By a linear constraint on the independent probability distribution, we mean the constraint of the type*

$$\underline{b} \leq \sum_{i=1}^{n_1} a_i \cdot p_i + \sum_{j=1}^{n_2} b_j \cdot q_j + \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} c_{i,j} \cdot p_i \cdot q_j \leq \bar{b}$$

for some given values \underline{b} , \bar{b} , a_i , b_j , and $c_{i,j}$.

- *By an interval probability problem under independence constraint, we mean the following problem: given a finite list of linear constraints, check whether there exists an independent distribution which satisfies all these constraints.*

Comment. Independence means that $p_{i,j} = p_i \cdot q_j$ for every i and j . The above constraints are linear in terms of these probabilities $p_{i,j} = p_i \cdot q_j$.

Theorem 3 *Interval probability problem under independence constraint is NP-hard.*

Proof. To prove this theorem, we will reduce the problem in question to the same known NP-hard problem as in the proof of Theorem 1: to the partition problem.

For every instance of the partition problem, we form the following system of constraints: $n_1 = n_2 = n$,

- $p_i - q_i = 0$ for every i from 1 to n ;
- $S_i \cdot p_i - p_i \cdot q_i = 0$ for all i from 1 to n ,

where

$$S_i \stackrel{\text{def}}{=} \frac{2 \cdot s_i}{\sum_{k=1}^n s_k}.$$

Let us prove that this system is consistent if and only if the original instance of the partition problem has a solution.

Indeed, if the original instance has a solution $\varepsilon_i \in \{-1, 1\}$, then, for every i from 1 to n , we can take $p_i = q_i = \frac{1 + \varepsilon_i}{2} \cdot S_i$, i.e.:

- if $\varepsilon_i = -1$, we take $p_i = q_i = 0$;
- if $\varepsilon_i = 1$, we take $p_i = q_i = S_i$.

Let us show that for this choice, $\sum_{i=1}^n p_i = \sum_{j=1}^n q_j = 1$. Indeed,

$$\sum_{i=1}^n p_i = \sum_{i=1}^n \frac{1 + \varepsilon_i}{2} \cdot S_i = \frac{1}{2} \cdot \sum_{i=1}^n S_i + \frac{1}{2} \cdot \sum_{i=1}^n \varepsilon_i \cdot S_i.$$

By definition of $S_i = \frac{2 \cdot s_i}{\sum_{k=1}^n s_k}$, we have

$$\sum_{i=1}^n S_i = 2 \cdot \frac{\sum_{i=1}^n s_i}{\sum_{k=1}^n s_k} = 2,$$

and

$$\sum_{i=1}^n \varepsilon_i \cdot S_i = 2 \cdot \frac{\sum_{i=1}^n \varepsilon_i \cdot s_i}{\sum_{k=1}^n s_k}.$$

Since $\sum_{i=1}^n \varepsilon_i \cdot s_i = 0$, the second sum is 0, hence $\sum_{i=1}^n p_i = 1$.

In both cases $\varepsilon_i = \pm 1$, we have $S_i \cdot p_i - p_i \cdot q_i = 0$, so all the constraints are indeed satisfied.

Vice versa, if the constraints are satisfied, this means that for every i , we have $p_i = q_i$ and $S_i \cdot p_i - p_i \cdot q_i = p_i \cdot (S_i - q_i) = p_i \cdot (S_i - p_i) = 0$, so $p_i = 0$ or $p_i = S_i$. Thus, the value p_i/S_i is equal to 0 or 1, hence the value $\varepsilon_i \stackrel{\text{def}}{=} 2 \cdot (p_i/S_i) - 1$ takes values -1 or 1 . In terms of ε_i , we have $p_i/S_i = \frac{1 + \varepsilon_i}{2}$, hence $p_i = \frac{1 + \varepsilon_i}{2} \cdot S_i$.

Since $\sum_{i=1}^n p_i = 1$, we conclude that

$$\sum_{i=1}^n p_i = \frac{1}{2} \cdot \sum_{i=1}^n S_i + \frac{1}{2} \cdot \sum_{i=1}^n \varepsilon_i \cdot S_i = 1.$$

We know that $\frac{1}{2} \cdot \sum_{i=1}^n S_i = 1$, hence $\sum_{i=1}^n \varepsilon_i \cdot S_i = 0$. We know that this sum is proportional to $\sum_{i=1}^n \varepsilon_i \cdot s_i$, hence $\sum_{i=1}^n \varepsilon_i \cdot s_i = 0$ – i.e., the original instance of the partition problem has a solution.

The theorem is proven.

Acknowledgments.

This work was supported in part by NASA under cooperative agreement NCC5-209, NSF grant EAR-0225670, NIH grant 3T34GM008048-20S1, and Army Research Lab grant DATM-05-02-C-0046.

The authors are very thankful to the participants of the 4th International Symposium on Imprecise Probabilities and Their Applications ISIPTA'05 (Carnegie Mellon University, July 20–24, 2005), especially to Mikelis Bickis (University of Saskatchewan, Canada), Arthur P. Dempster (Harvard University), and Damjan Škulj (University of Ljubljana, Slovenia), for valuable discussions.

References

- [1] D. Berleant, M.-P. Cheong, C. Chu, Y. Guan, A. Kamal, G. Sheblé, S. Ferson, and J. F. Peters, Dependable handling of uncertainty, *Reliable Computing*, 2003, Vol. 9, No. 6, pp. 407–418.
- [2] D. Berleant, M. Dancre, J. Argaud, and G. Sheblé, Electric company portfolio optimization under interval stochastic dominance constraints, In: F. G. Cozman, R. Nau, and T. Seidenfeld, *Proceedings of the 4th International Symposium on Imprecise Probabilities and Their Applications ISIPTA'05*, Pittsburgh, Pennsylvania, July 20–24, 2005, pp. 51–57.
- [3] D. J. Berleant, O. Kosheleva, and H. T. Nguyen, “Adding Unimodality or Independence Makes Interval Probability Problems NP-Hard”, *Proceedings of the International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems IPMU'06*, Paris, France, July 2–7, 2006 (to appear).
- [4] D. Berleant, L. Xie, and J. Zhang, Statool: a tool for Distribution Envelope Determination (DEnv), an interval-based algorithm for arithmetic on random variables, *Reliable Computing*, 2003, Vol. 9, No. 2, pp. 91–108.
- [5] D. Berleant and J. Zhang, Using Pearson correlation to improve envelopes around the distributions of functions, *Reliable Computing*, 2004, Vol. 10, No. 2, pp. 139–161.
- [6] D. Berleant and J. Zhang, Representation and Problem Solving with the Distribution Envelope Determination (DEnv) Method, *Reliability Engineering and System Safety*, 2004, Vol., 85, No. 1–3.
- [7] D. Berleant and J. Zhang, Using Pearson correlation to improve envelopes around the distributions of functions, *Reliable Computing*, 2004, Vol. 10, No. 2, pp. 139–161.
- [8] A. Borglin and H. Keiding, Stochastic dominance and conditional expectation an insurance theoretical approach, *The Geneva Papers on Risk and Insurance Theory*, 2002, Vol. 27, pp. 31-48.
- [9] A. Chateauneuf, On the use of capacities in modeling uncertainty aversion and risk aversion, *Journal of Mathematical Economics*, 1991, Vol. 20, pp. 343-369.
- [10] S. Ferson, *RAMAS Risk Calc 4.0*, CRC Press, Boca Raton, Florida.
- [11] M. R. Garey and D. S. Johnson, *Computers and Intractability, a Guide to the Theory of NP-Completeness*, W.H. Freeman and Company, San Francisco, CA, 1979.

- [12] G. J. Klir, G. Xiang, and O. Kosheleva, “Estimating information amount under interval uncertainty: algorithmic solvability and computational complexity”, *Proceedings of the International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems IPMU’06*, Paris, France, July 2–7, 2006 (to appear).
- [13] V. Kreinovich and S. Ferson, Computing Best-Possible Bounds for the Distribution of a Sum of Several Variables is NP-Hard, *International Journal of Approximate Reasoning*, 2006, Vol. 41, pp. 331–342.
- [14] V. Kreinovich, A. Lakeyev, J. Rohn, P. Kahl, *Computational complexity and feasibility of data processing and interval computations*, Kluwer, Dordrecht, 1997.
- [15] V. G. Krymsky, Computing Interval Estimates for Components of Statistical Information with Respect to Judgements on Probability Density Functions, In: J. Dongarra, K. Madsen, and J. Wasniewski (eds.), *PARA’04 Workshop on State-of-the-Art in Scientific Computing*, Springer Lecture Notes in Computer Science, 2005, Vol. 3732, pp. 151–160.
- [16] V. Kuznetsov, *Interval Statistical Models*, Radio i Svyaz, Moscow, 1991 (in Russian).
- [17] V. P. Kuznetsov, Interval methods for processing statistical characteristics, *Proceedings of the International Workshop on Applications of Interval Computations APIC’95*, El Paso, Texas, February 23–25, 1995 (a special supplement to the journal *Reliable Computing*), pp. 116–122.
- [18] V. P. Kuznetsov, Auxiliary problems of statistical data processing: interval approach, *Proceedings of the International Workshop on Applications of Interval Computations APIC’95*, El Paso, Texas, February 23–25, 1995 (a special supplement to the journal *Reliable Computing*), pp. 123–129.
- [19] D. Skulj, “Generalized conditioning in neighbourhood models”, In: F. G. Cozman, R. Nau, and T. Seidenfeld, *Proceedings of the 4th International Symposium on Imprecise Probabilities and Their Applications ISIPTA’05*, Pittsburgh, Pennsylvania, July 20–24, 2005.
- [20] D. Skulj, *A role of Jeffrey’s rule of conditioning in neighborhood models*, to appear.
- [21] P. Walley, *Statistical Reasoning with Imprecise Probabilities*, Chapman & Hall, N.Y., 1991.
- [22] J. Zhang and D. Berleant, Arithmetic on random variables: squeezing the envelopes with new joint distribution constraints, *Proceedings of the 4th International Symposium on Imprecise Probabilities and Their Applications ISIPTA’05*, Pittsburgh, Pennsylvania, July 20–24, 2005, pp. 416–422.