

A randomized method for handling a difficult function in a convex optimization problem, motivated by probabilistic programming

Csaba I. Fábíán*

Tamás Szántai†

Abstract

We propose a randomized gradient method for the handling of a convex function whose gradient computation is demanding. The method bears a resemblance to the stochastic approximation family. But in contrast to stochastic approximation, the present method builds a model problem.

The approach requires that estimates of function values and gradients be provided at the iterates. We present a variance reduction Monte Carlo simulation procedure to provide such estimates in the case of certain probabilistic functions.

Keywords. Convex optimization, stochastic optimization, probabilistic problems.

1 Introduction

We deal with approximate methods for the solution of smooth convex programming problems. First we consider minimization over a polyhedron:

$$\min \phi(T\mathbf{x}) \quad \text{subject to} \quad A\mathbf{x} \leq \mathbf{b}, \quad (1)$$

where $\phi : \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex function whose gradient computation is demanding. The vectors are $\mathbf{x} \in \mathbb{R}^m$, $\mathbf{b} \in \mathbb{R}^r$, and the matrices T and A are of sizes $n \times m$ and $r \times m$, respectively. For the sake of simplicity we assume that the feasible domain is not empty and is bounded. We then consider the minimization of a linear cost function subject to a difficult convex constraint:

$$\min \mathbf{c}^T \mathbf{x} \quad \text{subject to} \quad \check{A}\mathbf{x} \leq \check{\mathbf{b}}, \quad \phi(T\mathbf{x}) \leq \pi, \quad (2)$$

where the vectors $\mathbf{c}, \check{\mathbf{b}}$ and the matrix \check{A} have compatible sizes, and π is a given number.

The motivation for the above forms are the classic probability maximization and probabilistic constrained problems, where $\phi(\mathbf{z}) = -\log F(\mathbf{z})$ with a logconcave distribution function $F(\mathbf{z})$. In [7], an inner approximation was proposed for the probabilistic function. The approach proved easy to implement and invulnerable to noise in gradient computation. – Noisy gradient estimates may yield iterates that do not improve much on our current model. But we retain a true inner approximation, provided objective values at new iterates are evaluated with appropriate accuracy.

Let us briefly overview a couple of closely related probabilistic programming approaches. – For a broader survey, see [7] and references therein. – Given a distribution and a number p ($0 < p < 1$), a probabilistic constraint confines search to the level set $\mathcal{L}(F, p) = \{\mathbf{z} \mid F(\mathbf{z}) \geq p\}$ of the distribution function $F(\mathbf{z})$. Prékopa

*Department of Informatics, Faculty of Engineering and Computer Science, Kecskemét College, Pallasz Athéné University, Izsáki út 10, 6000 Kecskemét, Hungary. Email: fabian.csaba@gamf.kefo.hu.

†Department of Differential Equations, Faculty of Natural Sciences, Budapest University of Technology and Economics. Műegyetem rkp 3-9, 1111 Budapest, Hungary. Email: szantai@math.bme.hu.

[18] initiated a novel solution approach by introducing the concept of p-efficient points. \mathbf{z} is p-efficient if $F(\mathbf{z}) \geq p$ and there exists no \mathbf{z}' such that $\mathbf{z}' \leq \mathbf{z}$, $\mathbf{z}' \neq \mathbf{z}$, $F(\mathbf{z}') \geq p$. Prékopa, Vizvári, and Badics [20] consider problems with random parameters having a discrete finite distribution. They first enumerate p-efficient points, and based on these, build a convex relaxation of the problem.

Dentcheva, Prékopa, and Ruszczyński [6] formulate the probabilistic constraint in a split form: $T\mathbf{x} = \mathbf{z}$, where \mathbf{z} belongs to the level set $\mathcal{L}(F, p)$; and construct a Lagrangian dual by relaxing the constraint $T\mathbf{x} = \mathbf{z}$. The dual functional is the sum of two functionals that are respective optimal objective value functions of two simpler problems. The first auxiliary problem is a linear programming problem, and the second one is the minimization of a linear function over the level set $\mathcal{L}(F, p)$. Based on this decomposition, the authors develop a method, called cone generation, that finds new p-efficient points in course of the optimization process.

[7] focusses on probability maximization. A polyhedral approximation is constructed to the epigraph of the probabilistic function. This is analogous to the use of p-efficient points, and the approach was motivated by that concept. The dual problem is constructed and decomposed in the manner of [6], but the nonlinear subproblem is easier. In [6], finding a new p-efficient point amounts to minimization over the level set $\mathcal{L}(F, p)$. In contrast, a new approximation point in [7] is found by unconstrained minimization. Moreover, a practical approximation scheme was developed in [7]: instead of exactly solving an unconstrained subproblem occurring during the process, just a single line search was made in each case. Implementation based on this approximation scheme proved quite robust, and a theoretic explanation for this behavior was also found.

In the present paper, we extend the approach of [7] to handling gradient estimates. We also propose simulation schemes to obtain such estimates in case of probabilistic problems. We present the inner approximation approach in an idealized setting:

Assumption 1 *The function $\phi(\mathbf{z})$ is twice continuously differentiable, and the Hessian matrix satisfies*

$$\alpha I \preceq \nabla^2 \phi(\mathbf{z}) \preceq \omega I \quad (\mathbf{z} \in \mathbb{R}^n)$$

with some $\alpha, \omega \in \mathbb{R}$ ($0 < \alpha \leq \omega$). Here I is the identity matrix, and the relation $U \preceq V$ between matrices means that $V - U$ is positive semidefinite.

In Sections 2 and 3 we present a brisk overview of the models and the column generation approach of [7]. In Section 4 we present the column-generation approach from a dual viewpoint, as a cutting-plane method. (The dual viewpoint has the advantage that the cutting-plane method can be regularized, but we do not consider regularization in this paper.)

The cutting-plane model of the dual approach – like the inner approximation of the primal one – is invulnerable to gradient computation errors. This feature facilitates the use of gradient estimates. In Section 5 we extend the method in this direction. The motivation for applying gradient estimates was our computational experience reported in [7]: most of the computation effort was spent in computing gradients. – In that computational study we solved classic probability maximization problems; namely, we had $\phi(\mathbf{z}) = -\log F(\mathbf{z})$ with a multivariate normal distribution function $F(\mathbf{z})$. Given an n -dimensional normal distribution, a component of the gradient $\nabla F(\mathbf{z})$ was obtained as the product of an appropriate $(n - 1)$ -dimensional normal distribution function value and a univariate normal density function value (see the formula in Section 6.6.4 of Prékopa’s book [19]). The numerical computation of multivariate normal distribution function values was performed by Genz’s subroutine [11]. In our study, most of the computation time was spent in the Genz subroutine. Most demanding were the gradient computations, each requiring n calls to the Genz subroutine. – We conclude that easily computable estimates for the gradients are well worth using, even if the iteration count increases due to estimation errors.

For the estimation of function values and gradients in case of a probabilistic objective, we present a variance reduction Monte Carlo simulation procedure in Section 6. This procedure is applicable to gradient estimation in case of normal, Dirichlet, and t-distributions.

In Section 7 we deal with the convex constrained problem (2). Well-known approximation schemes for this problem consist of the solution of a sequence of problems of the form (1). We are going to show that the approximation tools described in previous sections facilitate such a solution scheme in the present case.

2 Problem and model formulation

In this section we formulate the dual problem and construct polyhedral models of the primal and dual problems. We follow the construction in [7], details can be found in that paper. Though in [7], we exploited monotonicity of the probabilistic objective and variable splitting was based on $\mathbf{z} \leq T\mathbf{x}$. In the present paper, we apply the traditional form of variable splitting: Problem (1) will be written as

$$\min \phi(\mathbf{z}) \quad \text{subject to} \quad A\mathbf{x} - \mathbf{b} \leq \mathbf{0}, \quad \mathbf{z} - T\mathbf{x} = \mathbf{0}. \quad (3)$$

This problem has an optimal solution due to our assumption that on feasible domain of (1). Introducing the multiplier vector $-\mathbf{y} \in \mathbb{R}^r$, $-\mathbf{y} \geq \mathbf{0}$ to the constraint $A\mathbf{x} - \mathbf{b} \leq \mathbf{0}$, and $-\mathbf{u} \in \mathbb{R}^n$ to the constraint $\mathbf{z} - T\mathbf{x} = \mathbf{0}$, the Lagrangian dual of (3) can be written as

$$\max \{\mathbf{y}^T \mathbf{b} - \phi^*(\mathbf{u})\} \quad \text{subject to} \quad (\mathbf{y}, \mathbf{u}) \in \mathcal{D}, \quad (4)$$

where

$$\mathcal{D} := \{ (\mathbf{y}, \mathbf{u}) \in \mathbb{R}^{r+n} \mid \mathbf{y} \leq \mathbf{0}, \quad T^T \mathbf{u} = A^T \mathbf{y} \}. \quad (5)$$

According to the theory of convex duality, this problem has an optimal solution.

2.1 Polyhedral models

Suppose we have evaluated the function $\phi(\mathbf{z})$ at points \mathbf{z}_i ($i = 0, 1, \dots, k$); let us introduce the notation $\phi_i = \phi(\mathbf{z}_i)$ for respective objective values. An inner approximation of $\phi(\cdot)$ is

$$\phi_k(\mathbf{z}) = \min \sum_{i=0}^k \lambda_i \phi_i \quad \text{such that} \quad \lambda_i \geq 0 \quad (i = 0, \dots, k), \quad \sum_{i=0}^k \lambda_i = 1, \quad \sum_{i=0}^k \lambda_i \mathbf{z}_i = \mathbf{z}. \quad (6)$$

If $\mathbf{z} \notin \text{Conv}(\mathbf{z}_0, \dots, \mathbf{z}_k)$, then let $\phi_k(\mathbf{z}) := +\infty$. A polyhedral model of Problem (3) is

$$\min \phi_k(\mathbf{z}) \quad \text{subject to} \quad A\mathbf{x} - \mathbf{b} \leq \mathbf{0}, \quad \mathbf{z} - T\mathbf{x} = \mathbf{0}. \quad (7)$$

We assume that (7) is feasible, i.e., its optimum is finite. This can be ensured by proper selection of the initial $\mathbf{z}_0, \dots, \mathbf{z}_k$ points. The convex conjugate of $\phi_k(\mathbf{z})$ is

$$\phi_k^*(\mathbf{u}) = \max_{0 \leq i \leq k} \{\mathbf{u}^T \mathbf{z}_i - \phi_i\}. \quad (8)$$

As $\phi_k^*(\cdot)$ is a cutting-plane model of $\phi^*(\cdot)$, the following problem is a polyhedral model of Problem (4):

$$\max \{\mathbf{y}^T \mathbf{b} - \phi_k^*(\mathbf{u})\} \quad \text{subject to} \quad (\mathbf{y}, \mathbf{u}) \in \mathcal{D}. \quad (9)$$

2.2 Linear programming formulations

The primal model problem (6)-(7) will be formulated as

$$\begin{aligned} \min \quad & \sum_{i=0}^k \phi_i \lambda_i \\ \text{such that} \quad & \lambda_i \geq 0 \quad (i = 0, \dots, k), \\ & \sum_{i=0}^k \lambda_i = 1, \\ & \sum_{i=0}^k \lambda_i \mathbf{z}_i - T\mathbf{x} = \mathbf{0}, \\ & A\mathbf{x} \leq \mathbf{b}. \end{aligned} \quad (10)$$

The dual model problem (8)-(9), formulated as a linear programming problem, is just the LP dual of (10):

$$\begin{aligned}
\max \quad & \vartheta \quad + \quad \mathbf{b}^T \mathbf{y} \\
\text{such that} \quad & \mathbf{y} \leq \mathbf{0}, \\
& \vartheta + \mathbf{z}_i^T \mathbf{u} \leq \phi_i \quad (i = 0, \dots, k), \\
& -T^T \mathbf{u} + A^T \mathbf{y} = \mathbf{0}.
\end{aligned} \tag{11}$$

Let $(\bar{\lambda}_0, \dots, \bar{\lambda}_k, \bar{\mathbf{x}})$ and $(\bar{\vartheta}, \bar{\mathbf{u}}, \bar{\mathbf{y}})$ denote respective optimal solutions of the problems (10) and (11) – both existing due to our assumption concerning the feasibility of (7) and hence (10). Let moreover

$$\bar{\mathbf{z}} = \sum_{i=0}^k \bar{\lambda}_i \mathbf{z}_i. \tag{12}$$

Observation 2 We have $\phi_k(\bar{\mathbf{z}}) = \sum_{i=0}^k \phi_i \bar{\lambda}_i = \bar{\vartheta} + \bar{\mathbf{u}}^T \bar{\mathbf{z}}$.

The first equality follows from the equivalence of (10) on the one hand, and (6)-(7) on the other hand. The second equality is a straight consequence of complementarity.

Observation 3 We have $\bar{\vartheta} = -\phi_k^*(\bar{\mathbf{u}})$.

This follows from the equivalence between (11) on the one hand and (8)-(9) on the other hand.

Remark 4 A consequence of Observations 2 and 3 is $\phi_k(\bar{\mathbf{z}}) + \phi_k^*(\bar{\mathbf{u}}) = \bar{\mathbf{u}}^T \bar{\mathbf{z}}$. This is Fenchel's equality between $\bar{\mathbf{u}}$ and $\bar{\mathbf{z}}$, with respect to the model function $\phi_k(\cdot)$.

3 Primal viewpoint: column generation

In [7], the probability maximization problem is solved by iteratively adding improving columns to the primal model. In this section we give a brisk overview of the practical approximation scheme proposed in that paper.

An optimal dual solution (i.e., shadow price vector) of the current model problem is $(\bar{\vartheta}, \bar{\mathbf{u}}, \bar{\mathbf{y}})$. Given a vector $\mathbf{z} \in \mathbb{R}^n$, we can add a new column in (10), corresponding to $\mathbf{z}_{k+1} = \mathbf{z}$. This is an improving column if its reduced cost

$$\bar{\rho}(\mathbf{z}) := \bar{\vartheta} + \bar{\mathbf{u}}^T \mathbf{z} - \phi(\mathbf{z}) \tag{13}$$

is positive. – It is easily seen that the reduced cost of $\bar{\mathbf{z}}$ is non-negative. Indeed,

$$\bar{\rho}(\bar{\mathbf{z}}) \geq \bar{\vartheta} + \bar{\mathbf{u}}^T \bar{\mathbf{z}} - \phi_k(\bar{\mathbf{z}}) = 0 \tag{14}$$

follows from $\phi_k(\cdot) \geq \phi(\cdot)$ and Observation 2.

In the context of the simplex method, the Markowitz rule is a well-known and often-used rule of column selection. The Markowitz rule selects the vector with the largest reduced cost. Coming back to the present problem (10), let

$$\bar{\mathcal{R}} := \max_{\mathbf{z}} \bar{\rho}(\mathbf{z}). \tag{15}$$

The column with the largest reduced cost can, in theory, be found by a steepest descent method applied to the function $-\bar{\rho}(\mathbf{z})$. Though finding a near-optimal solution proved rather time-consuming in the computational study of [7]. As a practical alternative, only a single line search was performed, starting from $\bar{\mathbf{z}}$. This simple method proved effective and robust. Moreover, a theoretical explanation was also found for the efficiency of the approach, based on the following well-known theorem:

Theorem 5 *Let Assumption 1 hold for the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Let us minimize $f(\mathbf{z})$ over \mathbb{R}^n using a steepest descent method, starting from a point \mathbf{z}^0 . Let $\mathbf{z}^1, \dots, \mathbf{z}^j, \dots$ denote the iterates obtained by applying exact line search at each step. Then we have*

$$f(\mathbf{z}^j) - \mathcal{F} \leq \left(1 - \frac{\alpha}{\omega}\right)^j [f(\mathbf{z}^0) - \mathcal{F}], \quad (16)$$

where $\mathcal{F} = \min_{\mathbf{z}} f(\mathbf{z})$.

Proof of this theorem can be found in, e.g., [22], [15]. The following corollary was obtained in [7]:

Corollary 6 *Let β ($0 < \beta \ll 1$) be given. A finite (and moderate) number of steps with the steepest descent method results a vector $\hat{\mathbf{z}}$ such that*

$$\bar{\rho}(\hat{\mathbf{z}}) \geq (1 - \beta) \bar{\mathcal{R}}. \quad (17)$$

This can be proven by substituting $f(\mathbf{z}) = -\bar{\rho}(\mathbf{z})$, $\mathbf{z}^0 = \bar{\mathbf{z}}$ in (16), and applying (14). Selecting j such that $(1 - \frac{\alpha}{\omega})^j \leq \beta$ yields an appropriate $\hat{\mathbf{z}} = \mathbf{z}^j$.

In view of the Markowitz rule mentioned above, the vector $\hat{\mathbf{z}}$ in Corollary 6 is a fairly good improving vector in the column generation scheme.

In the computational study of [7], just a single line search was performed in each reduced cost maximization; i.e., $j = 1$ according to the notation of Theorem 5. (Even this single line search was inexact, making a limited number of steps in the direction of steepest ascent.) Our implementation proved reliable even with this simple procedure.

In case of probabilistic functions, Assumption 1 does not hold for every $\mathbf{z} \in \mathbb{R}^n$. Our computational experience in [7] was, however, that the probabilistic objectives were well conditioned over certain domains. The iterates obtained by the above approximation procedure always remained in the respective safe domains.

Remark 7 *To check near-optimality of the current solution, we can use the usual LP stopping rule: $\bar{\mathcal{R}}$ should be less than a fixed optimality tolerance. – For the present special linear programming problem (10), this is not just a heuristic rule; as we show in the next section, $\bar{\mathcal{R}}$ is actually a bound on the gap between the respective optima of the model problem and the original convex problem.*

If we have good estimates for α and ω in Assumption 1, then $\bar{\mathcal{R}}$ can be estimated on the basis of Corollary 6. Hence the gap can be kept under effective control in course of the solution process.

If no reliable estimates for α and ω are known, then Corollary 6 is just a theoretical justification for limiting the numbers of the line searches in each steepest ascent method. The column generation procedure is terminated if $\|\nabla \bar{\rho}(\bar{\mathbf{z}})\|$ is small. We can construct an upper bound of the final gap $\bar{\mathcal{R}}$ using the gradient.

4 Dual viewpoint: cutting planes

The simplex method can be viewed as a cutting-plane method. This fact has been part of the professional folklore ever since the simplex method became widely known. Simplex and cutting-plane methods are parallelly discussed in Section 3.4 of Prékopa's book [19]. A closer description of the present situation can be found in [9], Section 4.

4.1 Dimension reduction

For the sake of simplicity, let us make

Assumption 8 *The inequality system $A\mathbf{x} \leq \mathbf{b}$ contains box constraints in the form of $\underline{\mathbf{b}} \leq \mathbf{x} \leq \bar{\mathbf{b}}$, where $\underline{\mathbf{b}}, \bar{\mathbf{b}} \in \mathbb{R}^r$ are given vectors ($\underline{\mathbf{b}} \leq \bar{\mathbf{b}}$).*

I.e., we have

$$A = \begin{pmatrix} A' \\ I \\ -I \end{pmatrix} \quad \text{and} \quad \mathbf{b} = \begin{pmatrix} \mathbf{b}' \\ \bar{\mathbf{b}} \\ -\bar{\mathbf{b}} \end{pmatrix}. \quad (18)$$

Let

$$\mu(\mathbf{u}) = \min \{ -\mathbf{b}^T \mathbf{y} \mid (\mathbf{y}, \mathbf{u}) \in \mathcal{D} \}. \quad (19)$$

This function is defined for every \mathbf{u} since $T^T \mathbf{u} = A^T \mathbf{y}$ is solvable in \mathbf{y} ($\mathbf{y} \leq \mathbf{0}$) due to A having the form (18). We will formulate the dual problems using the function $\mu(\mathbf{u})$. For convenience, we transform the problems into minimization forms. The original dual problem (4) assumes the form

$$\min \{ \phi^*(\mathbf{u}) + \mu(\mathbf{u}) \}, \quad (20)$$

and the model problem (9) assumes the form

$$\min \{ \phi_k^*(\mathbf{u}) + \mu(\mathbf{u}) \}. \quad (21)$$

4.2 Cut generation

Let $\bar{\mathbf{u}}$ be a minimizer of (21), in accordance with our notation in former sections. We are going to compute an approximate support function to $\phi^*(\mathbf{u})$ at $\bar{\mathbf{u}}$. This will be of the form

$$\ell(\mathbf{u}) := \mathbf{u}^T \mathbf{z} - \phi(\mathbf{z}), \quad (22)$$

with an appropriate vector \mathbf{z} . We have $\ell(\mathbf{u}) \leq \phi^*(\mathbf{u})$ for any \mathbf{u} by the definition of $\phi^*(\mathbf{u})$. We are going to compute \mathbf{z} such that $\phi^*(\bar{\mathbf{u}}) - \ell(\bar{\mathbf{u}})$ be relatively small.

Using support functions of the above form, a cutting-plane scheme for the problem (20) is easily implemented. We build a polyhedral model $\phi_k^*(\mathbf{u})$ of $\phi^*(\mathbf{u})$ – always adding the appropriate \mathbf{z} vector to the dual model as \mathbf{z}_{k+1} . On the other hand, we will work with $\mu(\mathbf{u})$ as a polyhedral model of itself. This is a workable setup; the current model function can be minimized by just solving the linear programming problem (11) – take into account Observation 3.

Coming back to the construction of the approximate support function (22), we wish to construct $\ell(\mathbf{u})$ whose graph cuts deeply into the epigraph of the model function $\phi_k^*(\mathbf{u})$. (Depth being measured at $\bar{\mathbf{u}}$.) I.e., we wish the following difference to be large:

$$\ell(\bar{\mathbf{u}}) - \phi_k^*(\bar{\mathbf{u}}) = \bar{\mathbf{u}}^T \mathbf{z} - \phi(\mathbf{z}) - \phi_k^*(\bar{\mathbf{u}}) = \bar{\rho}(\mathbf{z}), \quad (23)$$

where the second equality follows from Observation 3. (This is a dual interpretation the reduced cost.)

Let Ψ^* denote the minimum of (20). Let us further introduce the notation

$$\underline{\Psi} := \phi_k^*(\bar{\mathbf{u}}) + \mu(\bar{\mathbf{u}}) \quad \text{and} \quad \bar{\Psi} := \phi^*(\bar{\mathbf{u}}) + \mu(\bar{\mathbf{u}}). \quad (24)$$

Obviously we have $\underline{\Psi} \leq \Psi^* \leq \bar{\Psi}$. As for the gap between the upper and the lower bound, we have

$$\bar{\Psi} - \underline{\Psi} = \phi^*(\bar{\mathbf{u}}) - \phi_k^*(\bar{\mathbf{u}}) = \max_{\mathbf{z}} \bar{\rho}(\mathbf{z}) = \bar{\mathcal{R}}, \quad (25)$$

where the second equality follows from the definition of the conjugate function, and the third equality is in accordance with our former notation.

In order to construct inexact cuts, let us consider a dual form of Corollary 6:

Proposition 9 *Let β ($0 < \beta \ll 1$) be given. We can construct a linear function $\hat{\ell}(\mathbf{u})$ such that*

$$\hat{\ell}(\mathbf{u}) \leq \phi^*(\mathbf{u}) \quad \text{holds for any } \mathbf{u}, \text{ and}$$

$$\hat{\ell}(\bar{\mathbf{u}}) \geq \phi^*(\bar{\mathbf{u}}) - \beta \bar{\mathcal{R}}.$$

In words: $\hat{\ell}(\mathbf{u})$ is an approximate support function to $\phi^(\mathbf{u})$ at $\bar{\mathbf{u}}$. The difference between the function values at the current iterate is bounded by the portion $\beta \bar{\mathcal{R}}$ of the gap.*

Proof. According to Corollary 6, we can construct a vector $\widehat{\mathbf{z}}$ such that (17) holds. Using this $\widehat{\mathbf{z}}$, let us define $\hat{\ell}(\mathbf{u})$ according to (22). Then we have

$$\begin{aligned}\hat{\ell}(\bar{\mathbf{u}}) &= \bar{\rho}(\widehat{\mathbf{z}}) + \phi_k^*(\bar{\mathbf{u}}) && \text{due to (23)} \\ &\geq (1 - \beta)\bar{\mathcal{R}} + \phi_k^*(\bar{\mathbf{u}}) && \text{due to (17)} \\ &= \phi^*(\bar{\mathbf{u}}) - \beta\bar{\mathcal{R}} && \text{due to (25)}.\end{aligned}\tag{26}$$

5 Working with gradient estimates

In this section we show that the column generation scheme of [7] (sketched in Section 3), and the cutting-plane scheme of Section 4 can be implemented as a randomized method using gradient estimates.

We wish to minimize $-\bar{\rho}(\mathbf{z})$ over \mathbb{R}^n . Given $\mathbf{z}^\circ \in \mathbb{R}^n$, let $\mathbf{g}^\circ = -\nabla\bar{\rho}(\mathbf{z}^\circ)$.

Assumption 10 *Let $\sigma > 0$ be given. We can construct a realization of a random vector \mathbf{G}° satisfying*

$$E(\mathbf{G}^\circ) = \mathbf{g}^\circ \quad \text{and} \quad E\left(\|\mathbf{G}^\circ - \mathbf{g}^\circ\|^2\right) \leq \sigma \|\mathbf{g}^\circ\|^2.\tag{27}$$

From (27) follows

$$E\left(\|\mathbf{G}^\circ\|^2\right) = E\left(\|\mathbf{G}^\circ - \mathbf{g}^\circ\|^2\right) + \|\mathbf{g}^\circ\|^2 \leq (\sigma + 1) \|\mathbf{g}^\circ\|^2.\tag{28}$$

Let us consider the following randomized form of Theorem 5:

Theorem 11 *Let Assumption 1 hold for the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Let us minimize $f(\mathbf{z})$ over \mathbb{R}^n . We perform a steepest descent method using gradient estimates. (Given an iterate \mathbf{z}° , a gradient estimate \mathbf{G}° is generated and a line search is performed in that direction.) We assume that gradient estimates at the respective iterates are generated independently, and (27) - (28) hold for each of them.*

Having started from the point \mathbf{z}^0 , and having performed j line searches, let $\mathbf{z}^1, \dots, \mathbf{z}^j$ denote the respective iterates. Then we have

$$E[f(\mathbf{z}^j)] - \mathcal{F} \leq \left(1 - \frac{\alpha}{\omega(\sigma + 1)}\right)^j (f(\mathbf{z}^0) - \mathcal{F}),\tag{29}$$

where $\mathcal{F} = \min_{\mathbf{z}} f(\mathbf{z})$.

Proof. Let $\mathbf{G}^0, \dots, \mathbf{G}^{j-1}$ denote the respective gradient estimates for the iterates $\mathbf{z}^0, \dots, \mathbf{z}^{j-1}$.

To begin with, let us focus on the first line search whose starting point is $\mathbf{z}^\circ = \mathbf{z}^0$. Here \mathbf{z}° is a given (not random) vector. We are going to adopt the usual proof of Theorem 5 to employing the gradient estimate \mathbf{G}° instead of the gradient \mathbf{g}° . From $\nabla^2 f(\mathbf{z}) \preceq \omega I$, it follows that

$$f(\mathbf{z}^\circ - t\mathbf{G}^\circ) \leq f(\mathbf{z}^\circ) - t\mathbf{g}^{\circ T}\mathbf{G}^\circ + \frac{\omega}{2}t^2\mathbf{G}^{\circ T}\mathbf{G}^\circ$$

holds for any $t \in \mathbb{R}$. Considering expectations in both sides, we get

$$\begin{aligned}E[f(\mathbf{z}^\circ - t\mathbf{G}^\circ)] &\leq f(\mathbf{z}^\circ) - t\|\mathbf{g}^\circ\|^2 + \frac{\omega}{2}t^2 E\left(\|\mathbf{G}^\circ\|^2\right) \\ &\leq f(\mathbf{z}^\circ) - t\|\mathbf{g}^\circ\|^2 + \frac{\omega}{2}t^2(\sigma + 1)\|\mathbf{g}^\circ\|^2\end{aligned}$$

according to (28). Let us consider the respective minima in t separately of the two sides. The right-hand side is a quadratic expression, yielding minimum at $t = \frac{1}{\omega(\sigma + 1)}$. Inequality is inherited to minima, hence

$$\min_t E[f(\mathbf{z}^\circ - t\mathbf{G}^\circ)] \leq f(\mathbf{z}^\circ) - \frac{1}{2\omega(\sigma + 1)}\|\mathbf{g}^\circ\|^2.\tag{30}$$

For the left-hand side, we obviously have

$$\mathbb{E} \left[\min_t f(\mathbf{z}^\circ - t\mathbf{G}^\circ) \right] \leq \min_t \mathbb{E} [f(\mathbf{z}^\circ - t\mathbf{G}^\circ)]. \quad (31)$$

(This is analogous to the basic inequality comparing the wait-and-see and the here-and-now approaches for classic two-stage stochastic programming problems.)

Let \mathbf{z}' denote the minimizer of the line search in the left-hand side of (31), i.e., $f(\mathbf{z}') = \min_t f(\mathbf{z}^\circ - t\mathbf{G}^\circ)$. (Of course \mathbf{z}' is a random vector since it depends on \mathbf{G}° .) Substituting this in (31) and comparing with (30), we get

$$\mathbb{E} [f(\mathbf{z}')] \leq f(\mathbf{z}^\circ) - \frac{1}{2\omega(\sigma+1)} \|\mathbf{g}^\circ\|^2.$$

Subtracting \mathcal{F} from both sides results

$$\mathbb{E} [f(\mathbf{z}')] - \mathcal{F} \leq f(\mathbf{z}^\circ) - \mathcal{F} - \frac{1}{2\omega(\sigma+1)} \|\mathbf{g}^\circ\|^2. \quad (32)$$

Coming to the lower bound, a well-known consequence of $\alpha I \preceq \nabla^2 f(\mathbf{z})$ is

$$\|\mathbf{g}^\circ\|^2 \geq 2\alpha (f(\mathbf{z}^\circ) - \mathcal{F})$$

(see the classic proof of Theorem 5). Combining this with (32), we get

$$\mathbb{E} [f(\mathbf{z}')] - \mathcal{F} \leq f(\mathbf{z}^\circ) - \mathcal{F} - \frac{\alpha}{\omega(\sigma+1)} (f(\mathbf{z}^\circ) - \mathcal{F}) = \left(1 - \frac{\alpha}{\omega(\sigma+1)}\right) (f(\mathbf{z}^\circ) - \mathcal{F}). \quad (33)$$

As we have assumed that \mathbf{z}° is a given (not random) vector, the right-hand side of (33) is deterministic, and the expectation in the left-hand side is considered according to the distribution of \mathbf{G}° .

Let us now examine the $(l+1)$ th line search (for $1 \leq l \leq j-1$) where the starting point is $\mathbf{z}^\circ = \mathbf{z}^l$ and the minimizer is $\mathbf{z}' = \mathbf{z}^{l+1}$. Of course (33) holds with these objects also, but now both sides are random variables, depending on the vectors $\mathbf{G}^0, \dots, \mathbf{G}^{l-1}$. (The expectation in the left-hand side is a conditional expectation.) Let us consider the respective expectations of the two sides, according to the joint distribution of $\mathbf{G}^0, \dots, \mathbf{G}^{l-1}$. As the random gradient vectors were generated independently, we get

$$\mathbb{E} [f(\mathbf{z}^{l+1})] - \mathcal{F} \leq \left(1 - \frac{\alpha}{\omega(\sigma+1)}\right) (\mathbb{E} [f(\mathbf{z}^l)] - \mathcal{F}), \quad (34)$$

where the left-hand expectation is now taken according to the joint distribution of $\mathbf{G}^0, \dots, \mathbf{G}^l$. – This technique of proof is well-known in the context of stochastic gradient schemes.

Finally, (29) follows from the iterative application of (34). \square

Corollary 12 *Let a tolerance β ($0 < \beta \ll 1$) and a probability p ($0 < p \ll 1$) be given. A finite (and moderate) number of steps with the above randomized steepest descent method results a vector $\hat{\mathbf{z}}$ such that*

$$P(\bar{\rho}(\hat{\mathbf{z}}) \geq (1 - \beta)\bar{\mathcal{R}}) \geq 1 - p.$$

I.e., with a high probability, $\hat{\mathbf{z}}$ is a fairly good improving vector in the column generation scheme.

Proof. Substituting $f(\mathbf{z}) = -\bar{\rho}(\mathbf{z})$ and $\mathbf{z}^0 = \bar{\mathbf{z}}$ in (29) and taking into account (14), we get

$$\mathbb{E} [\bar{\rho}(\mathbf{z}^j)] \geq (1 - \varrho^j) \bar{\mathcal{R}}$$

with $\varrho = 1 - \frac{\alpha}{\omega(\sigma+1)}$. The gap $\bar{\mathcal{R}}$ is obviously non-negative. In case $\bar{\mathcal{R}} = 0$, the starting iterate $\mathbf{z}^0 = \bar{\mathbf{z}}$ of the steepest descent method was already optimal, due to (14). In what follows we assume $\bar{\mathcal{R}} > 0$. A trivial transformation results

$$\mathbb{E} \left[1 - \frac{\bar{\rho}(\mathbf{z}^j)}{\bar{\mathcal{R}}} \right] \leq \varrho^j.$$

By Markov's inequality, we get

$$\mathbb{P}\left(1 - \frac{\bar{\rho}(\mathbf{z}^j)}{\bar{\mathcal{R}}} \geq \beta\right) \leq \frac{\varrho^j}{\beta},$$

and a trivial transformation yields

$$\mathbb{P}\left(\bar{\rho}(\mathbf{z}^j) \leq (1 - \beta)\bar{\mathcal{R}}\right) \leq \frac{1}{\beta} \varrho^j.$$

Hence

$$\mathbb{P}\left(\bar{\rho}(\mathbf{z}^j) > (1 - \beta)\bar{\mathcal{R}}\right) \geq 1 - \frac{1}{\beta} \varrho^j.$$

Selecting j such that $\varrho^j \leq \beta p$ yields an appropriate $\hat{\mathbf{z}} = \mathbf{z}^j$. \square

Gradients of the function $-\bar{\rho}(\mathbf{z})$ have the form $\nabla\phi(\mathbf{z}) - \bar{\mathbf{u}}$. As the procedure progresses, the difference $\nabla\phi(\hat{\mathbf{z}}) - \bar{\mathbf{u}}$ gets small. To satisfy the requirement (27) on variance, better and better estimates are needed. We can control accuracy like in the deterministic column generation method of Section 3.

Remark 13 *If we have good estimates for α and ω in Assumption 1, then $\bar{\mathcal{R}}$ can be estimated on the basis of Corollary 12. Otherwise Corollary 12 is just a theoretical justification for limiting the numbers of the line searches in each steepest ascent procedure.*

When the column generation scheme stops, we need to check the magnitude of the current gap. It may happen that only a statistical verification is possible, on the basis of Assumption 10.

A dual form of the above corollary is

Proposition 14 *Let a tolerance β ($0 < \beta \ll 1$) and a probability p ($0 < p \ll 1$) be given. We can construct a random linear function $\hat{\ell}(\mathbf{u})$ such that*

$$\hat{\ell}(\mathbf{u}) \leq \phi^*(\mathbf{u}) \quad \text{holds for any } \mathbf{u}, \text{ and}$$

$$\mathbb{P}\left(\hat{\ell}(\bar{\mathbf{u}}) \geq \phi^*(\bar{\mathbf{u}}) - \beta\bar{\mathcal{R}}\right) \geq 1 - p.$$

The proof is the same as that of Proposition 9, but use Corollary 12 instead of Corollary 6.

6 Easily computable estimates of function values and gradients

For the partial derivatives of any multivariate probability distribution function we have the general formula

$$\frac{\partial F(z_1, \dots, z_n)}{\partial z_i} = F(z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_n | z_i) f_i(z_i) \quad (35)$$

where $F(z_1, \dots, z_n)$ is the probability distribution function of the random variables ξ_1, \dots, ξ_n , $F(z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_n | z_i)$ is the conditional probability distribution function of the random variables $\xi_1, \dots, \xi_{i-1}, \xi_{i+1}, \dots, \xi_n$, given that $\xi_i = z_i$ and $f_i(z)$ is the probability density function of the random variable ξ_i (see Formula (6.6.22) on the page 203 of Prékopa's book [19]). From Formula (35) it follows that if the multivariate probability distribution at issue has conditional probability distributions of its own type then we can calculate the multivariate probability distribution function values and its partial derivatives by the same procedure. Such type of multivariate probability distributions are for example the multivariate normal, the multivariate t -distribution and the Dirichlet distribution.

Therefore in this section we present a variance reduction Monte Carlo simulation procedure for the estimation of multivariate probability distribution function values, only. The procedure was proposed in Szántai's thesis [23], and quoted in Sections 6.5 and 6.6 of Prékopa's book [19].

This procedure can be applied to any multivariate probability distribution function. The only condition is that we have to be able to calculate the one- and the two-dimensional marginal probability distribution function values. Accuracy can easily be controlled by changing the sample size.

This way we can construct gradient estimates satisfying Assumption 10.

As we have

$$F(z_1, \dots, z_n) = P(\xi_1 < z_1, \dots, \xi_n < z_n) = 1 - P(\bar{A}_1 \cup \dots \cup \bar{A}_n),$$

where

$$\bar{A}_i = \{\xi_i \geq z_i\} \quad (i = 1, \dots, n),$$

we can apply bounding and simulation results for the probability of union of events.

6.1 Crude Monte Carlo simulation

If μ denotes the number of those events which occur out of the events $\bar{A}_1, \bar{A}_2, \dots, \bar{A}_n$, then the random variable

$$\nu_0 = \begin{cases} 0, & \text{if } \mu = 0 \\ 1, & \text{if } \mu \geq 1 \end{cases}$$

obviously has expected value $\bar{P} = P(\bar{A}_1 \cup \bar{A}_2 \cup \dots \cup \bar{A}_n)$.

6.2 Monte Carlo simulation of the differences between the true probability value and its respective second order Boole–Bonferroni bounds

For the probability $\bar{P} = P(\bar{A}_1 \cup \bar{A}_2 \cup \dots \cup \bar{A}_n)$ we have the so called second order Boole–Bonferroni bounds (see Prékopa's book [19]).

The lower bound is

$$L_2 = \frac{2}{k^* + 1} \bar{S}_1 - \frac{2}{k^*(k^* + 1)} \bar{S}_2,$$

where

$$k^* = 1 + \left\lfloor \frac{2\bar{S}_2}{\bar{S}_1} \right\rfloor.$$

The upper bound is

$$U_2 = \bar{S}_1 - \frac{2}{n} \bar{S}_2.$$

In both cases \bar{S}_1 and \bar{S}_2 are the first resp. second order binomial moments of the random variable μ which can be expressed also as

$$\bar{S}_1 = \sum_{i=1}^n P(\bar{A}_i), \quad \bar{S}_2 = \sum_{1 \leq i_1 < i_2 \leq n} P(\bar{A}_{i_1} \cap \bar{A}_{i_2}).$$

Then by applying the Poincare (sieve) formula we get for the differences between the probability value \bar{P} and its second order Boole–Bonferroni bounds

$$\bar{P} - L_2 = \bar{S}_1 - \bar{S}_2 + \dots + (-1)^{n-1} \bar{S}_n - \frac{2}{k^* + 1} \bar{S}_1 + \frac{2}{k^*(k^* + 1)} \bar{S}_2,$$

and

$$\bar{P} - U_2 = -\bar{S}_2 + \dots + (-1)^{n-1} \bar{S}_n + \frac{2}{n} \bar{S}_2.$$

So the random variables

$$\nu_{L_2} = \begin{cases} 0, & \text{if } \mu \leq 1 \\ \sum_{i=1}^n (-1)^{i-1} \binom{\mu}{i} - \frac{2}{k^*+1} \binom{\mu}{1} + \frac{2}{k^*(k^*+1)} \binom{\mu}{2} = \frac{1}{k^*(k^*+1)} (\mu - k^*)(\mu - k^* - 1), & \text{if } \mu \geq 2 \end{cases}$$

and

$$\nu_{U_2} = \begin{cases} 0, & \text{if } \mu \leq 1 \\ \sum_{i=2}^n (-1)^{i-1} \binom{\mu}{i} - \frac{2}{n} \binom{\mu}{2} = \frac{1}{n} (\mu+n)(1-\mu), & \text{if } \mu \geq 2 \end{cases}$$

have expected values $\bar{P} - L_2$ and $\bar{P} - U_2$. This way the transformed random variables $\nu_{L_2} + L_2$ and $\nu_{U_2} + U_2$ also have expected value \bar{P} .

6.3 Monte Carlo simulation of the difference between the true probability value and its Hunter–Worsley bound

The Hunter–Worsley upper bound for the probability $\bar{P} = P(\bar{A}_1 \cup \bar{A}_2 \cup \dots \cup \bar{A}_n)$ is defined as (see Prékopa's book [19]):

$$U_{HW} = \bar{S}_1 - \sum_{(i,j) \in T^*} P(\bar{A}_i \cap \bar{A}_j) \geq \bar{P},$$

where T^* is the maximum weight spanning tree in the complete graph with nodes $\{1, \dots, n\}$ and edges $\{\{i, j\}, 1 \leq i, j \leq n\}$ and to node i the probability $P(\bar{A}_i)$ and to edge $\{i, j\}$ the probability $P(\bar{A}_i \cap \bar{A}_j)$ is assigned.

Then by applying the Poincare formula we get for the difference between the probability value \bar{P} and its Hunter–Worsley upper bound

$$\bar{P} - U_{HW} = -\bar{S}_2 + \bar{S}_3 - \dots + (-1)^{n-1} \bar{S}_n + \sum_{(i,j) \in T^*} P(\bar{A}_i \cap \bar{A}_j).$$

If λ denotes the number of those $\bar{A}_i \cap \bar{A}_j$, $(i, j) \in T^*$ events which occur in a random trial, then the random variable

$$\nu_{HW} = \begin{cases} 0, & \text{if } \mu \leq 1 \\ \sum_{i=2}^n (-1)^{i-1} \binom{\mu}{i} + \lambda = 1 - \mu + \lambda, & \text{if } \mu \geq 2 \end{cases}$$

has expected value $\bar{P} - U_{HW}$ and so the transformed random variable $\nu_{HW} + U_{HW}$ also has expected value \bar{P} .

6.4 Determination of the final estimation with minimal variance

Let us chose the random variables ν_0 , $\nu_{HW} + U_{HW}$ and $\nu_{L_2} + L_2$ and denote $\widehat{P}_0, \widehat{P}_1, \widehat{P}_2$ the three different estimations based on these. All of these are unbiased estimations of the probability value \bar{P} . Let us estimate the empirical covariances of these estimations in a simulation procedure:

$$\widehat{C} = \begin{pmatrix} \widehat{c}_{00} & \widehat{c}_{01} & \widehat{c}_{02} \\ \widehat{c}_{01} & \widehat{c}_{11} & \widehat{c}_{12} \\ \widehat{c}_{02} & \widehat{c}_{12} & \widehat{c}_{22} \end{pmatrix}$$

If we introduce the new estimation

$$\widehat{P} = w_0 \widehat{P}_0 + w_1 \widehat{P}_1 + w_2 \widehat{P}_2$$

where $w_0 + w_1 + w_2 = 1$, then it will be also an unbiased estimation of the probability value \bar{P} . As \widehat{P} has variance $w^T \widehat{C} w$, where $w^T = (w_0, w_1, w_2)$, therefore the coefficients w_0, w_1, w_2 resulting the minimal variance estimation can be determined by the solution of the nonlinear programming problem:

$$\min_{w_0+w_1+w_2=1} w^T \widehat{C} w.$$

As the gradient of $w^T \widehat{C} w$ equals to $2w^T \widehat{C}$ it is easy to see that the unknown values of w_1, w_2, w_3, λ can be determined by the solution of the following linear equation system:

$$\begin{aligned} \widehat{c}_{00}w_0 + \widehat{c}_{01}w_1 + \widehat{c}_{02}w_2 - \lambda &= 0, \\ \widehat{c}_{01}w_0 + \widehat{c}_{11}w_1 + \widehat{c}_{12}w_2 - \lambda &= 0, \\ \widehat{c}_{02}w_0 + \widehat{c}_{12}w_1 + \widehat{c}_{22}w_2 - \lambda &= 0, \\ w_0 + w_1 + w_2 &= 1, \end{aligned} \tag{36}$$

representing the Karush-Kuhn-Tucker necessary condition.

6.5 Further Monte Carlo simulation algorithms

For the case of multivariate normal probability distribution there are other known Monte Carlo simulation algorithms, see Deák [3], [4] and Ambartzumian et al [1]. Gassmann [10] combined Szántai's general algorithm and Deák's algorithm into a hybrid algorithm. The efficiency of this algorithm was explored by Deák, Gassmann and Szántai in [5].

7 Handling a difficult constraint

We are going to work out an approximation scheme for the solution of the convex constrained problem (2). This scheme will consist of the solution of a sequence of problems of the form (1).

Let us consider the linear constraint set $A\mathbf{x} \leq \mathbf{b}$ of Problem (1). The last constraint of this set is $\mathbf{a}^r \mathbf{x} \leq b_r$, where \mathbf{a}^r denotes the r th row of A , and b_r denotes the r th component of \mathbf{b} . Assume that this last constraint is a cost constraint, and let $\mathbf{c}^T = \mathbf{a}^r$ denote the cost vector. We are going to consider a parametric form of the cost constraint, namely, $\mathbf{c}^T \mathbf{x} \leq d$, where $d \in \mathbb{R}$ is a parameter.

Let \check{A} denote the matrix obtained by omitting the r th row in A , and let $\check{\mathbf{b}}$ denote the vector obtained by omitting the r th component in \mathbf{b} . Using these objects, let us consider the problem

$$\min \phi(T\mathbf{x}) \quad \text{subject to} \quad \check{A}\mathbf{x} \leq \check{\mathbf{b}}, \quad \mathbf{c}^T \mathbf{x} \leq d, \tag{37}$$

with the parameter $d \in \mathbb{R}$. This parametric form of the unconstrained problem will be denoted by (1: $b_r = d$).

Let $\chi(d)$ denote the optimal objective value of Problem (37), as a function of the parameter d . This is obviously a monotone decreasing convex function. Let $\mathcal{I} \subset \mathbb{R}$ denote the domain over which the function is finite. We have either $\mathcal{I} = \mathbb{R}$ or $\mathcal{I} = [\underline{d}, +\infty)$ with some $\underline{d} \in \mathbb{R}$. Using the notation of the unconstrained problem, we can say that $\chi(d)$ is the optimum of (1: $b_r = d$) for $d \in \mathcal{I}$.

Coming to the constrained problem (2), we assume that the right-hand value π has been set by an expert, on the basis of preliminary experimental information. We may assume $\pi \in \chi(\mathcal{I})$. Let $d^* \in \mathcal{I}$ be a solution of the equation $\chi(d) = \pi$, and let $l^*(d)$ denote a linear support function to $\chi(d)$ at d^* . Every decision maker would set a right-hand side π that satisfies

Assumption 15 *The support function $l^*(d)$ has a significant negative slope, i.e., $l^{*\prime} \ll 0$.*

It follows that the optimal objective value of (2) is d^* .

We are going to find a near-optimal $\widehat{d} \in \mathcal{I}$ using an approximate version of Newton's method. – The idea of regulating tolerances in such a procedure occurs in the discussion of the Constrained Newton Method in

[14]. Based on the convergence proof of the Constrained Newton Method, a simple convergence proof of Newton's method was reconstructed in [8]. We are going to adapt the latter to the present case.

A sequence of unconstrained problems will be solved with increasing accuracy. In course of this procedure, we are going to build a single model $\phi_k(\mathbf{z})$ of the nonlinear objective $\phi(\mathbf{z})$. (I.e., k is ever increasing.) Given an iterate $d^\circ \in \mathcal{I}$, we can estimate $\chi(d^\circ)$ by solving the current model problem (7: $b_r = d^\circ$). Using the notation of the dual approach, let the optimum of the current model problem be $-\underline{\Psi} = -\phi_k^*(\bar{u}) - \mu(\bar{u})$. According to (25), we have $-\underline{\Psi} \geq \chi(d^\circ) \geq -\underline{\Psi} - \bar{\mathcal{R}}$. If we can compute gradients of $\phi(\mathbf{z})$ then we can construct an upper bound of the gap $\bar{\mathcal{R}}$, as mentioned in Remark 7. If we work with gradient estimates then only a statistical estimation of final gap is possible, as mentioned in Remark 13. First we describe the approximation scheme for the former case. In order to handle gradient estimates, we work out a randomized version of the approximation scheme in Section 7.2.

7.1 A deterministic approximation scheme

In this section we consider the case when upper bounds can be constructed for the gap in the unconstrained problem. Let $d_0, d_1 \in \mathcal{I}$, $d_0 < d_1 < d^*$ be the starting iterates. – The sequence of the iterates will be strictly monotone increasing, and converging to d^* from below.

Near-optimality condition

Given a tolerance ϵ ($\pi \gg \epsilon > 0$), let $\hat{d} \in \mathcal{I}$ be such that

$$\hat{d} \leq d^* \quad \text{and} \quad \chi(\hat{d}) \leq \pi + \epsilon. \quad (38)$$

Let $\hat{\mathbf{x}}$ be an optimal solution of (37: $d = \hat{d}$). Then $\hat{\mathbf{x}}$ is an ϵ -feasible solution of (2) with objective value \hat{d} . Exact feasible solutions of (2) have objective values not less than $d^* \geq \hat{d}$.

Evaluation of $\chi(d_j)$

Given iterate $d_j \in \mathcal{I}$, $d_j \leq d^*$, we are going to include the unknown value $\chi(d_j)$ in a known interval whose length is comparable to $\chi(d_j) - \pi$. Namely, we are going to find an upper bound $\bar{\chi}_j$ such that

$$\bar{\chi}_j \geq \chi(d_j) \geq \bar{\chi}_j - \delta(\bar{\chi}_j - \pi), \quad (39)$$

where δ ($0 < \delta \ll \frac{1}{2}$) is a fixed tolerance.

Such $\bar{\chi}_j$ is found by approximately solving the problem (37: $d = d_j$) \equiv (1: $b_r = d_j$). Let us use the column generation / cutting plane schemes described in previous sections. The optimum of the current model problem is always an upper bound for $\chi(d_j)$. Let us stop the column generation / cutting plane procedure when either of the following condition is satisfied:

$$\begin{aligned} (i) \quad & \bar{\chi}_j - \pi \leq \epsilon, \quad \text{or} \\ (ii) \quad & \bar{\mathcal{R}}_j \leq \delta(\bar{\chi}_j - \pi), \end{aligned} \quad (40)$$

where $\bar{\chi}_j$ will be an upper approximation of $\chi(d_j)$, and $\bar{\mathcal{R}}_j$ will measure the quality of this approximation.

Namely, let $\bar{\chi}_j$ be the optimum of the current model problem (7: $b_r = d_j$). We have

$$\bar{\chi}_j \geq \chi(d_j). \quad (41)$$

Let moreover $\bar{\mathcal{R}}_j$ be an upper estimate of the gap in the current model, i.e., of $\bar{\mathcal{R}}$. We have

$$\bar{\chi}_j - \chi(d_j) \leq \bar{\mathcal{R}}_j. \quad (42)$$

The tools described in Sections 3 and 4 allow us, in principle, to decrease $\bar{\mathcal{R}}_j$ below any tolerance.

If (i) occurs then, taking into account (41), $\hat{d} := d_j$ satisfies the near-optimality condition (38).

If (ii) occurs then, taking into account (41)-(42), $\bar{\chi}_j$ satisfies (39), yielding a usable approximation of $\chi(d_j)$.

Finding successive iterates

Given $j \geq 1$, assume that we have evaluated $\chi(d_{j-1})$ and $\chi(d_j)$, as described in the previous section. We are going to find d_{j+1} . In order to do this, let us first consider the linear function $l_j(d)$ satisfying

$$l_j(d_{j-1}) := \bar{\chi}_{j-1} \geq \chi(d_{j-1}) \quad \text{and} \quad l_j(d_j) := \bar{\chi}_j - \delta(\bar{\chi}_j - \pi) \leq \chi(d_j), \quad (43)$$

where the inequalities follow from (39). Due to the convexity of $\chi(d)$ and to Assumption 15, the above linear function obviously has a negative slope $l'_j \leq l^{\star'} \ll 0$. Moreover $l_j(d) \leq \chi(d)$ holds for $d_j \leq d$.

The next iterate d_{j+1} will be the point where $l_j(d_{j+1}) = \pi$. Of course $d_j < d_{j+1} \leq d^*$ follows from the observations above.

Convergence

Let the iterates d_0, d_1, \dots, d_s and the linear functions $l_1(d), \dots, l_s(d)$ be as defined above. We assume that $s > 1$, and the procedure did not stop before step $(s + 1)$. Then we have

$$\bar{\chi}_j - \pi > \epsilon \quad (j = 0, 1, \dots, s). \quad (44)$$

To simplify notation, let us introduce $L_j(d) = l_j(d) - \pi$ ($j = 1, \dots, s$). We transform (43) into

$$L_j(d_{j-1}) = \bar{\chi}_{j-1} - \pi \quad \text{and} \quad L_j(d_j) = (1 - \delta)(\bar{\chi}_j - \pi) \quad (j = 1, \dots, s), \quad (45)$$

positivity of function values following from (44). Moreover, the derivatives satisfy

$$L'_j = l'_j \leq l^{\star'} \ll 0 \quad (j = 1, \dots, s) \quad (46)$$

due to the observations in the previous section.

Theorem 16 *We have*

$$\gamma^{(s-1)} \mathcal{K} L_1(d_1) \geq L_s(d_s), \quad (47)$$

where $0 < \gamma \ll 1$ and \mathcal{K} is a positive number of moderate magnitude.

Proof. The following statements hold for $j = 1, \dots, s - 1$. From (45), we get

$$\frac{L_{j+1}(d_j)}{L_j(d_j)} = \frac{\bar{\chi}_j - \pi}{(1 - \delta)(\bar{\chi}_j - \pi)} = \frac{1}{1 - \delta}. \quad (48)$$

By definition, we have

$$L_j(d_j) + (d_{j+1} - d_j) L'_j = L_j(d_{j+1}) = 0.$$

It follows that $d_{j+1} - d_j = \frac{L_j(d_j)}{|L'_j|}$. Using this, we get

$$L_{j+1}(d_j) = L_{j+1}(d_{j+1}) + (d_j - d_{j+1}) L'_{j+1} = L_{j+1}(d_{j+1}) + \frac{L_j(d_j)}{|L'_j|} |L'_{j+1}|.$$

Hence

$$\frac{L_{j+1}(d_j)}{L_j(d_j)} = \frac{L_{j+1}(d_{j+1})}{L_j(d_j)} + \frac{|L'_{j+1}|}{|L'_j|}. \quad (49)$$

From (48), we have

$$\frac{1}{1 - \delta} = \frac{L_{j+1}(d_{j+1})}{L_j(d_j)} + \frac{|L'_{j+1}|}{|L'_j|} \geq 2 \sqrt{\frac{L_{j+1}(d_{j+1}) |L'_{j+1}|}{L_j(d_j) |L'_j|}}.$$

(This is the well-known inequality between means.) It follows that

$$\left(\frac{1}{2(1-\delta)}\right)^2 L_j(d_j) |L'_j| \geq L_{j+1}(d_{j+1}) |L'_{j+1}|.$$

By induction, we get

$$\left(\frac{1}{2(1-\delta)}\right)^{2(s-1)} L_1(d_1) |L'_1| \geq L_s(d_s) |L'_s|. \quad (50)$$

Dividing both sides by $|L'_s|$, and substituting

$$\gamma := \left(\frac{1}{2(1-\delta)}\right)^2, \quad \mathcal{K} := \frac{|L'_1|}{|L'^*|} \geq \frac{|L'_1|}{|L'_s|},$$

we obtain (47). Here $\gamma \ll 1$ follows from the setting of the constant $0 < \delta \ll \frac{1}{2}$. \mathcal{K} has a moderate magnitude according to (46), eventually due to Assumption 15.

7.2 A randomized version of the approximation scheme

In this section we consider the case when gradient estimates are used in solving the unconstrained problems. Only a statistical estimation of the gap $\overline{\mathcal{R}}$ is possible, on the basis of Assumption 10. Let $\overline{\mathcal{R}}_j$ denote our upper estimate, in accordance with the notation used for the deterministic scheme. We may underestimate the gap, meaning that (42) does not hold; and consequently the right-hand inequality of (39) may not hold. The probability of such an occurrence can be kept low. In such exceptional cases, $d_{j+1} > d^*$ may occur.

Hence we need to check $\chi(d_{j+1})$ in the new iterate. If we can verify $\chi(d_{j+1}) > \pi$ then we can proceed as in the deterministic scheme. If we can verify $\chi(d_{j+1}) < \pi$ then we can just step back to the previous iterate d_j and re-estimate the gap, possibly with a higher reliability. (The process may involve tightening of the upper bound $\overline{\chi}_j$.)

It may happen that neither $\chi(d_{j+1}) > \pi$ nor $\chi(d_{j+1}) < \pi$ can be verified because $\chi(d_{j+1})$ is near to π . The latter fact can also be verified, and in this case d_{j+1} can be considered a near-optimal solution to (2), under mild additional assumptions on $\chi(\cdot)$ and π .

8 Conclusion and discussion

In this paper we present the column-generation approach of [7] from a dual viewpoint, as a cutting-plane method. Moreover we propose a randomized version of this method. There is an important contrast between direct cutting-plane methods and the present approach. Direct cutting-plane methods for probabilistic functions are difficult to implement due to noise in gradient computation. Even a fairly accurate gradient may result a cut cutting into the epigraph of the objective function (especially in regions farther away from the current iterate). One either needs sophisticated tolerance handling to avoid cutting into the epigraph – see, e.g., [24], [16], [2] –, or else one needs a sophisticated convex optimization method that can handle cuts cutting into the epigraph – see [26]. (Yet another alternative, developed for a different type of problem, is perpetual adjustment of existing cuts to information revealed in course of the process; see [12].)

The present models are invulnerable to gradient computation errors. Noisy gradient estimates may yield iterates that do not improve much on our current models. But we retain a true inner approximation of the primal objective – or a true outer approximation of the dual objective –, provided objective values at new iterates are evaluated with appropriate accuracy. This feature facilitates the use of gradient estimates. Our randomized method bears a resemblance to the stochastic approximation family that goes back to [21] (see [17], [13] for recent forms).

The use of gradient estimates may substantially decrease total computational effort, even though a certain (moderate) accuracy *is* demanded in objective values. Computing a single component of a gradient

vector will involve an effort comparable to that of computing an objective value, e.g., in case of probability maximization under multivariate normal distribution of the random parameters.

The variance reduction Monte Carlo simulation procedure described in Section 6 was successfully applied in the solution of jointly probabilistic constrained stochastic programming problems, see [24]. The situation was similar to the present one; as the procedure progressed, higher and higher accuracy became necessary.

The approximation scheme of Section 7 consist of the solution of a sequence of problems of the form (1), i.e., minimization of a convex objective over polyhedra. Suppose the problems in this sequence are solved by the approximation approach described in previous sections. Then can we build a single model of the nonlinear objective $\phi(\mathbf{z})$. I.e., the model built in course of the solution of problem \mathcal{P}_s can be used as a starting model for the solution of the successive problem \mathcal{P}_{s+1} .

Future work

Our motivation for dealing with a difficult function was a probabilistic function $F(\mathbf{z}) = P(\boldsymbol{\xi} < \mathbf{z})$, where the random vector $\boldsymbol{\xi}$ has a logconcave distribution. The proposed approach can be extended to two-sided probabilities of the type

$$P(\underline{\mathbf{z}} < \boldsymbol{\xi} < \bar{\mathbf{z}}), \quad (51)$$

where $\underline{\mathbf{z}}$ and $\bar{\mathbf{z}}$ are linear functions of the decision variables, i.e., we have $\underline{\mathbf{z}} = \underline{T}\mathbf{x} + \underline{\mathbf{t}}$ and $\bar{\mathbf{z}} = \bar{T}\mathbf{x} + \bar{\mathbf{t}}$ with appropriate matrices \underline{T}, \bar{T} and vectors $\underline{\mathbf{t}}, \bar{\mathbf{t}}$. The Monte Carlo simulation procedures described in Section 6 can be applied by using

$$\bar{A}_i = \{\xi_i \leq z_i\} \cup \{\bar{z}_i \leq \xi_i\} \quad (i = 1, \dots, n).$$

Van Ackooij, Henrion, Möller and Zorgati [25] developed gradient formulas for two-sided probabilities (51) in case of normal distributions. Analogous formulas can be developed for other multivariate probability distributions.

Acknowledgement

This research is supported by the EFOP-3.6.1-16-2016-00006 "The development and enhancement of the research potential at Pallas Athena University" project. The Project is supported by the Hungarian Government and co-financed by the European Social Fund.

References

- [1] R. Ambartzumian, A. Der Kiureghian, V. Ohanian, and H. Sukiasian. Multinormal probability by sequential conditioned importance sampling: Theory and applications. *Probabilistic Engineering Mechanics*, 13:299–308, 1998.
- [2] T. Arnold, R. Henrion, A. Möller, and S. Vigerske. A mixed-integer stochastic nonlinear optimization problem with joint probabilistic constraints. *Pacific Journal of Optimization*, 10:5–20, 2014.
- [3] I. Deák. Three digit accurate multiple normal probabilities. *Numerische Mathematik*, 35:369–380, 1980.
- [4] I. Deák. Computing probabilities of rectangles in case of multinormal distributions. *Journal of Statistical Computation and Simulation*, 26:101–114, 1986.
- [5] I. Deák, H. Gassmann, and T. Szántai. Computing multivariate normal probabilities: a new look. *Journal of Statistical Computation and Simulation*, 11:920–949, 2002.
- [6] D. Dentcheva, A. Prékopa, and A. Ruszczyński. Concavity and efficient points of discrete distributions in probabilistic programming. *Mathematical Programming*, 89:55–77, 2000.

- [7] C.I. Fábián, E. Csizmás, R. Drenyovszki, W. van Ackooij, T. Vajnai, L. Kovács, and T. Szántai. Probability maximization by inner approximation. *Submitted to Acta Polytechnica Hungarica*, 2017.
- [8] C.I. Fábián, K. Eretnek, and O. Papp. A regularized simplex method. *Central European Journal of Operations Research*, 23:877–898, 2015.
- [9] C.I. Fábián, O. Papp, and K. Eretnek. Implementing the simplex method as a cutting-plane method, with a view to regularization. *Computational Optimization and Applications*, 56:343–368, 2013.
- [10] H. Gassmann. Conditional probability and conditional expectation of a random vector. In Y.M. Ermoliev and R.J.-B. Wets, editors, *Numerical Techniques for Stochastic Optimization*, pages 237–254. Springer-Verlag, Berlin, 1988.
- [11] A. Genz. Numerical computation of multivariate normal probabilities. *Journal of Computational and Graphical Statistics*, 1:141–150, 1992.
- [12] J.L. Higle and S. Sen. *Stochastic Decomposition: A Statistical Method for Large Scale Stochastic Linear Programming*, volume 8 of *Nonconvex Optimization and Its Applications*. Springer, 1996.
- [13] Guanghui Lan. An optimal method for stochastic composite optimization. *Mathematical Programming*, 133:365–397, 2012.
- [14] C. Lemaréchal, A. Nemirovskii, and Yu. Nesterov. New variants of bundle methods. *Mathematical Programming*, 69:111–147, 1995.
- [15] D.G. Luenberger and Y. Ye. *Linear and Nonlinear Programming*. International Series in Operations Research and Management Science. Springer, 2008.
- [16] J. Mayer. *Stochastic Linear Programming Algorithms: A Comparison Based on a Model Management System*. Gordon and Breach Science Publishers, 1998.
- [17] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19:1574–1609, 2009.
- [18] A. Prékopa. Dual method for a one-stage stochastic programming problem with random RHS obeying a discrete probability distribution. *ZOR - Methods and Models of Operations Research*, 34:441–461, 1990.
- [19] A. Prékopa. *Stochastic Programming*. Kluwer Academic Publishers, Dordrecht, 1995.
- [20] A. Prékopa, B. Vizvári, and T. Badics. Programming under probabilistic constraint with discrete random variable. In F. Giannesi, T. Rapsák, and S. Komlósi, editors, *New Trends in Mathematical Programming*, pages 235–255. Kluwer, Dordrecht, 1998.
- [21] H. Robbins and S. Monro. A stochastic approximation method. *Annals of Mathematical Statistics*, 22:400–407, 1951.
- [22] A. Ruszczyński. *Nonlinear Optimization*. Princeton University Press, 2006.
- [23] T. Szántai. *Numerical Evaluation of Probabilities Concerning Multidimensional Probability Distributions, Thesis*. Hungarian Academy of Sciences, Budapest, 1985.
- [24] T. Szántai. A computer code for solution of probabilistic-constrained stochastic programming problems. In Y.M. Ermoliev and R.J.-B. Wets, editors, *Numerical Techniques for Stochastic Optimization*, pages 229–235. Springer-Verlag, Berlin, 1988.
- [25] W. van Ackooij, R. Henrion, A. Möller, and R. Zorgati. On probabilistic constraints induced by rectangular sets and multivariate normal distributions. *Mathematical Methods of Operations Research*, 71:535–549, 2010.

- [26] W. van Ackooij and C. Sagastizábal. Constrained bundle methods for upper inexact oracles with application to joint chance constrained energy problems. *SIAM Journal on Optimization*, 24:733–765, 2014.