

A randomized method for smooth convex minimization, motivated by probability maximization

Csaba I. Fábíán*

Tamás Szántai[†]

Abstract

We propose a randomized gradient method – or a randomized cutting-plane method from a dual viewpoint. From the primal viewpoint, our method bears a resemblance to the stochastic approximation family. But in contrast to stochastic approximation, the present method builds a model problem.

Keywords. Convex optimization, stochastic optimization, probabilistic problems.

1 Introduction

We deal with approximate methods for the solution of the problem

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

where $\mathbf{x} \in \mathbb{R}^m$, $\mathbf{b} \in \mathbb{R}^r$, 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.

The motivation for the above form 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. This was motivated by the concept of p-efficient points introduced by Prékopa [16]. 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. – In this paper we present the inner approximation approach in an idealized setting:

Assumption 1 *The function $\phi : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex, twice continuously differentiable, and the Hessian matrix $\nabla^2\phi(\mathbf{z})$ satisfies*

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

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

We'll present a dual view of the method proposed in [7]. Then we'll propose a randomized version of that method: a randomized gradient method from a primal viewpoint – or a randomized cutting-plane method from a dual viewpoint. But first 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 [16] initiated a novel solution approach

*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.

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 [18] 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. 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.

The main result of the present paper is extending the approach of [7] to handling gradient estimates.

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 advantage of the dual viewpoint is that the cutting-plane method can be regularized. Actual means of regularization are not discussed in this paper, but we present the method in a form that admits of iterate selection rules other than plain cutting-plane method.

The cutting-plane model of the dual approach – like the inner approximation of the primal one – is invulnerable to gradient computation errors. 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. 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 [17]). The numerical computation of multivariate normal distribution function values was performed by Genz’s subroutine [10]. 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.

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

$$\min \phi(\mathbf{z}) \quad \text{subject to} \quad A\mathbf{x} - \mathbf{b} \leq \mathbf{0}, \quad \mathbf{z} - T\mathbf{x} = \mathbf{0}. \quad (2)$$

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 (2) can be written as

$$\max \{\mathbf{y}^T \mathbf{b} - \phi^*(\mathbf{u})\} \quad \text{subject to} \quad (\mathbf{y}, \mathbf{u}) \in \mathcal{D}, \quad (3)$$

where

$$\mathcal{D} := \{(\mathbf{y}, \mathbf{u}) \in \mathbb{R}^{r+n} \mid \mathbf{y} \leq \mathbf{0}, \quad T^T \mathbf{u} = A^T \mathbf{y}\}. \quad (4)$$

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

$$\phi_k(\mathbf{z}) = \min \sum_{i=0}^k \lambda_i \phi_i \quad \text{such that} \quad \lambda_i \geq 0 \quad (i = 0, \dots, k), \quad \sum_{i=0}^k \lambda_i = 1, \quad \sum_{i=0}^k \lambda_i \mathbf{z}_i = \mathbf{z}. \quad (5)$$

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 (2) is

$$\min \phi_k(\mathbf{z}) \quad \text{subject to} \quad A\mathbf{x} - \mathbf{b} \leq \mathbf{0}, \quad \mathbf{z} - T\mathbf{x} = \mathbf{0}. \quad (6)$$

We assume that (6) 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

$$\phi_k^*(\mathbf{u}) = \max_{0 \leq i \leq k} \{\mathbf{u}^T \mathbf{z}_i - \phi_i\}. \quad (7)$$

As $\phi_k^*(\cdot)$ is a cutting-plane model of $\phi^*(\cdot)$, the following problem is a polyhedral model of Problem (3):

$$\max \{\mathbf{y}^T \mathbf{b} - \phi_k^*(\mathbf{u})\} \quad \text{subject to} \quad (\mathbf{y}, \mathbf{u}) \in \mathcal{D}. \quad (8)$$

2.2 Linear programming formulations

The primal model problem (5)-(6) will be formulated as

$$\begin{aligned} \min \quad & \sum_{i=0}^k \phi_i \lambda_i \\ \text{such that} \quad & \lambda_i \geq 0 \quad (i = 0, \dots, k), \\ & \sum_{i=0}^k \lambda_i = 1, \\ & \sum_{i=0}^k \lambda_i \mathbf{z}_i - T\mathbf{x} = \mathbf{0}, \\ & A\mathbf{x} \leq \mathbf{b}. \end{aligned} \quad (9)$$

The dual model problem (7)-(8), formulated as a linear programming problem, is just the LP dual of (9):

$$\begin{aligned}
\max \quad & \vartheta \quad + \quad \mathbf{b}^T \mathbf{y} \\
\text{such that} \quad & \mathbf{y} \leq \mathbf{0}, \\
& \vartheta + \mathbf{z}_i^T \mathbf{u} \leq \phi_i \quad (i = 0, \dots, k), \\
& -T^T \mathbf{u} + A^T \mathbf{y} = \mathbf{0}.
\end{aligned} \tag{10}$$

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 (9) and (10) – both existing due to our assumption concerning the feasibility of (6) and hence (9). Let moreover

$$\bar{\mathbf{z}} = \sum_{i=0}^k \bar{\lambda}_i \mathbf{z}_i. \tag{11}$$

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 (9) on the one hand, and (5)-(6) 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 (10) on the one hand and (7)-(8) 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 (9), corresponding to $\mathbf{z}_{k+1} = \mathbf{z}$. This is an improving column if its reduced cost

$$\bar{\rho}(\mathbf{z}) := \bar{\vartheta} + \bar{\mathbf{u}}^T \mathbf{z} - \phi(\mathbf{z}) \tag{12}$$

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

$$\bar{\rho}(\bar{\mathbf{z}}) \geq \bar{\vartheta} + \bar{\mathbf{u}}^T \bar{\mathbf{z}} - \phi_k(\bar{\mathbf{z}}) = 0 \tag{13}$$

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

Let $\bar{\mathcal{R}} := \max_{\mathbf{z}} \bar{\rho}(\mathbf{z})$. We have $\bar{\mathcal{R}} = \bar{\vartheta} + \phi^*(\bar{\mathbf{u}})$ by the definition of the conjugate function. If $\bar{\mathcal{R}}$ is small, then $(\bar{\mathbf{x}}, \bar{\mathbf{z}})$ is a near-optimal solution to (2). Otherwise an improving column can be constructed to (9). 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*

$$f(\mathbf{z}^j) - \mathcal{F} \leq \left(1 - \frac{\alpha}{\omega}\right)^j [f(\mathbf{z}^0) - \mathcal{F}], \quad (14)$$

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

Proof of this theorem can be found in, e.g., [20], [13]. 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*

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

In the context of the simplex method, the Markowitz rule is a well-known and often-used rule of column selection. The Markowitz rule always selects the vector with the largest reduced cost. In this view, 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 approximative procedure always remained in the respective safe domains.

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 [17]. A closer description of the present situation can be found in [8], Section 4.

4.1 Dimension reduction

For the sake of simplicity, let us make

Assumption 7 *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

$$A = \begin{pmatrix} A' \\ I \\ -I \end{pmatrix} \quad \text{and} \quad \mathbf{b} = \begin{pmatrix} \mathbf{b}' \\ \bar{\mathbf{b}} \\ -\underline{\mathbf{b}} \end{pmatrix}. \quad (15)$$

Let

$$\mu(\mathbf{u}) = \min \{ -\mathbf{b}^T \mathbf{y} \mid (\mathbf{y}, \mathbf{u}) \in \mathcal{D} \}. \quad (16)$$

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 (15). 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 (3) assumes the form

$$\min \{ \phi^*(\mathbf{u}) + \mu(\mathbf{u}) \}, \quad (17)$$

and the model problem (8) assumes the form

$$\min \{ \phi_k^*(\mathbf{u}) + \mu(\mathbf{u}) \}. \quad (18)$$

4.2 Cut generation

Given $\hat{\mathbf{u}} \in \mathbb{R}^n$, we are going to compute an approximate support function to $\phi^*(\mathbf{u})$. This will be of the form

$$\hat{\ell}(\mathbf{u}) := \mathbf{u}^T \hat{\mathbf{z}} - \phi(\hat{\mathbf{z}}), \quad (19)$$

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

Using support functions of the above form, a cutting-plane scheme for the problem (17) is easily implemented. We build a polyhedral model $\phi_k^*(\mathbf{u})$ of $\phi^*(\mathbf{u})$ – always adding the appropriate $\hat{\mathbf{z}}$ 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 (10) – take into account Observation 3.

Let Ψ^* denote the minimum of (17). Moreover let $\bar{\mathbf{u}}$ be a minimizer of (18), in accordance with our notation in former sections. Let us further introduce the notation

$$\underline{\Psi} := \phi_k^*(\bar{\mathbf{u}}) + \mu(\bar{\mathbf{u}}) \quad \text{and} \quad \widehat{\Psi} := \phi^*(\hat{\mathbf{u}}) + \mu(\hat{\mathbf{u}}). \quad (20)$$

Obviously we have $\underline{\Psi} \leq \Psi^* \leq \widehat{\Psi}$.

Coming back to the construction of the support function (19), a deep cut at $\hat{\mathbf{u}}$ – relative to the model function $\phi_k^*(\mathbf{u})$ – can be obtained by setting $\hat{\mathbf{z}}$ to be a near-maximizer of

$$\widehat{\rho}(\mathbf{z}) = \hat{\mathbf{u}}^T \mathbf{z} - \phi(\mathbf{z}) - \phi_k^*(\hat{\mathbf{u}}). \quad (21)$$

Let $\widehat{\mathcal{R}} := \max_{\mathbf{z}} \widehat{\rho}(\mathbf{z})$. We have

$$\widehat{\mathcal{R}} = \phi^*(\hat{\mathbf{u}}) - \phi_k^*(\hat{\mathbf{u}}) \quad (22)$$

by the definition of the conjugate function. As for the gap between the upper and the lower bound, we have

$$\begin{aligned} \widehat{\Psi} - \underline{\Psi} &= \phi^*(\hat{\mathbf{u}}) - \phi_k^*(\hat{\mathbf{u}}) + \phi_k^*(\hat{\mathbf{u}}) - \phi_k^*(\bar{\mathbf{u}}) + (\mu(\hat{\mathbf{u}}) - \mu(\bar{\mathbf{u}})) \\ &= \widehat{\mathcal{R}} + (\phi_k^*(\hat{\mathbf{u}}) - \phi_k^*(\bar{\mathbf{u}})) + (\mu(\hat{\mathbf{u}}) - \mu(\bar{\mathbf{u}})), \end{aligned} \quad (23)$$

from (20) and (22). We are going to work with inexact cuts. Let us first consider a corollary of Theorem 5.

Corollary 8 *Let β ($0 < \beta \ll 1$) be given. A finite (and moderate) number of steps with the steepest descent method results a vector $\bar{\mathbf{z}}$ such that*

$$\widehat{\rho}(\bar{\mathbf{z}}) \geq (1 - \beta) \widehat{\mathcal{R}} + \beta \widehat{\rho}(\bar{\mathbf{z}}). \quad (24)$$

Proof. Substituting $f(\mathbf{z}) = -\widehat{\rho}(\mathbf{z})$ and $\mathbf{z}^0 = \bar{\mathbf{z}}$ in (14), and introducing the notation $\varrho = 1 - \alpha/\omega$, we get

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

(We have $\mathcal{F} = -\widehat{\mathcal{R}}$ by definition.) Selecting j such that $\varrho^j \leq \beta$ yields an appropriate $\hat{\mathbf{z}} = \mathbf{z}^j$. \square

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

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

$$\hat{\ell}(\hat{\mathbf{u}}) \geq \phi^*(\hat{\mathbf{u}}) - \beta \left[\widehat{\mathcal{R}} - \widehat{\rho}(\bar{\mathbf{z}}) \right].$$

Proof. According to Corollary 8, we can construct a vector $\widehat{\mathbf{z}}$ such that (24) holds. Using this $\widehat{\mathbf{z}}$, let us define $\widehat{\ell}(\mathbf{u})$ according to (19). Then we have

$$\begin{aligned}
\widehat{\ell}(\widehat{\mathbf{u}}) &= \widehat{\rho}(\widehat{\mathbf{z}}) + \phi_k^*(\widehat{\mathbf{u}}) && \text{due to (21)} \\
&\geq (1 - \beta) \widehat{\mathcal{R}} + \beta \widehat{\rho}(\widehat{\mathbf{z}}) + \phi_k^*(\widehat{\mathbf{u}}) && \text{due to (24)} \\
&= \phi^*(\widehat{\mathbf{u}}) - \beta \widehat{\mathcal{R}} + \beta \widehat{\rho}(\widehat{\mathbf{z}}) && \text{due to (22)}.
\end{aligned} \tag{25}$$

4.3 A special case

In the remaining part of this paper we consider a special case:

Assumption 10 *The next iterate $\widehat{\mathbf{u}}$ is always selected in such a way that*

$$\widehat{\rho}(\widehat{\mathbf{z}}) = \widehat{\mathbf{u}}^T \widehat{\mathbf{z}} - \phi(\widehat{\mathbf{z}}) - \phi_k^*(\widehat{\mathbf{u}}) \geq 0. \tag{26}$$

In words, (26) means that adding the support function (19: $\widehat{\mathbf{z}} = \widehat{\mathbf{z}}$) would improve the model function $\phi_k^*(\mathbf{u})$ at $\widehat{\mathbf{u}}$.

According to (13), this assumption holds in case of the plain cutting-plane method, where the next iterate is simply the minimizer of the model problem (18); namely, we have $\widehat{\mathbf{u}} = \bar{\mathbf{u}}$ and hence $\widehat{\rho}(\widehat{\mathbf{z}}) = \bar{\rho}(\widehat{\mathbf{z}})$.

Under Assumption 10, we can simplify Corollary 8 by deleting the term $\beta \widehat{\rho}(\widehat{\mathbf{z}})$ in the right-hand side of (24). In a similar manner, Proposition 9 will assume the form

Proposition 11 *Let β ($0 < \beta \ll 1$) be given. We can construct a linear function $\widehat{\ell}(\mathbf{u})$ such that*

$$\begin{aligned}
\widehat{\ell}(\mathbf{u}) &\leq \phi^*(\mathbf{u}) \quad \text{holds for any } \mathbf{u}, \text{ and} \\
\widehat{\ell}(\widehat{\mathbf{u}}) &\geq \phi^*(\widehat{\mathbf{u}}) - \beta \widehat{\mathcal{R}}.
\end{aligned}$$

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

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

Assumption 12 *Let $\sigma > 0$ be given. We can construct a realization of a random vector \mathbf{G}° satisfying*

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

From (27) follows

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

Let us consider the following randomized form of Theorem 5:

Theorem 13 *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

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

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

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

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

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

For the left-hand side, we obviously have

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

(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

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

Subtracting \mathcal{F} from both sides results

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

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

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

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

Under Assumption 10, we have

Corollary 14 *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*

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

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

Proof. Substituting $f(\mathbf{z}) = -\hat{\rho}(\mathbf{z})$ and $\mathbf{z}^0 = \bar{\mathbf{z}}$ in (29) and taking into account Assumption 10, we get

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

with $\varrho = 1 - \frac{\alpha}{\omega(\sigma+1)}$. We have $\hat{\mathcal{R}} \geq 0$ by (22). In case $\hat{\mathcal{R}} = 0$, the starting iterate $\mathbf{z}^0 = \bar{\mathbf{z}}$ of the steepest descent method was already optimal, due to Assumption 10. In what follows we assume $\hat{\mathcal{R}} > 0$. A trivial transformation results

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

By Markov's inequality, we get

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

and a trivial transformation yields

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

Hence

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

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

Remark 15 *We work with gradients of $f(\mathbf{z}) = -\hat{\rho}(\mathbf{z})$. By the definition of the latter function, we have*

$$\nabla f(\mathbf{z}) = \nabla \phi(\mathbf{z}) - \bar{\mathbf{u}}.$$

As the procedure progresses, the difference $\nabla \phi(\mathbf{z}^j) - \bar{\mathbf{u}}$ gets small. To satisfy the requirement (27) on variance, better and better estimates are needed.

5.1 Random cutting planes

Under Assumption 10 we can apply random cutting planes in the cutting-plane scheme of Section 4:

Proposition 16 *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*

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

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

The proof is the same as those of Propositions 9 and 11, but use Corollary 14 instead of Corollary 8.

6 Easily computable estimates of function values and gradients

In this section we present a variance reduction Monte Carlo simulation procedure for the estimation of multivariate probability distribution function values. The procedure was proposed in Szántai's thesis [21], and presented in Section 6.5 of Prékopa's book [17].

This procedure can be applied to any multivariate probability distribution function. The only condition is that we have to be able to calculate one- and two-dimensional marginal probability distribution function values. Accuracy is easily controlled by sample size. The procedure can be applied to evaluate objective values at new iterates with appropriate accuracy.

All these algorithms can be applied to gradient estimation if the multivariate probability distribution at issue has conditional probability distributions of its own type. E.g., the multivariate normal, the multivariate t -distribution and the Dirichlet distribution have this property. In such cases way we can effectively construct gradient estimates satisfying Assumption 12.

As we have

$$F(x_1, \dots, x_n) = P(\xi_1 < x_1, \dots, \xi_n < x_n) = 1 - P(\{\xi_1 \geq x_1\} + \dots + \{\xi_n \geq x_n\}) = 1 - P(\bar{A}_1 + \dots + \bar{A}_n),$$

where

$$\bar{A}_i = \{\xi_i \geq x_i\}, i = 1, \dots, n,$$

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 $\bar{P} = P(\bar{A}_1 + \bar{A}_2 + \dots + \bar{A}_n)$.

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

If we take L_2 as the second order Boole–Bonferroni lower bound of the probability of union of events and U_2 as the second order Boole–Bonferroni upper bound of the probability of union of events

$$L_2 = \frac{2}{k^* + 1} \bar{S}_1 - \frac{2}{k^*(k^* + 1)} \bar{S}_2 \leq \bar{P} \leq \bar{S}_1 - \frac{2}{n} \bar{S}_2 = U_2,$$

then

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

and

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

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

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

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 differences between the true probability value and its Hunter–Worsley bound

Let

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

then

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

If μ denotes the number of those events which occur out of the events $\bar{A}_1, \bar{A}_2, \dots, \bar{A}_n$, further λ denotes the number of those $\bar{A}_i \bar{A}_j$, $(i, j) \in T^*$ 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 \overline{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 Lagrangian 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:

$$\begin{aligned} \widehat{c}_{00}w_0 + \widehat{c}_{01}w_1 + \widehat{c}_{02}w_2 - \lambda &= 0, \\ \widehat{c}_{01}w_0 + \widehat{c}_{11}w_1 + \widehat{c}_{12}w_2 - \lambda &= 0, \\ \widehat{c}_{02}w_0 + \widehat{c}_{12}w_1 + \widehat{c}_{22}w_2 - \lambda &= 0, \\ w_0 + w_1 + w_2 &= 1. \end{aligned} \tag{35}$$

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 [9] 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 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., [22], [14], [2] –, or else one needs a sophisticated convex optimization method that can handle cuts cutting into the epigraph – see [23]. (Yet another alternative, developed for a different type of problem, is perpetual adjustment of the existing model to information revealed in course of the process; see [11].)

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 [19] (see [15], [12] 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 [22]. The situation was similar to the present one; as the procedure progressed, higher and higher accuracy became necessary.

Acknowledgement

This research is supported by 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, 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.
- [9] 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.
- [10] A. Genz. Numerical computation of multivariate normal probabilities. *Journal of Computational and Graphical Statistics*, 1:141–150, 1992.
- [11] J.L. Hight and S. Sen. *Stochastic Decomposition: A Statistical Method for Large Scale Stochastic Linear Programming*, volume 8 of *Nonconvex Optimization and Its Applications*. Springer, 1996.
- [12] Guanghui Lan. An optimal method for stochastic composite optimization. *Mathematical Programming*, 133:365–397, 2012.
- [13] D.G. Luenberger and Y. Ye. *Linear and Nonlinear Programming*. International Series in Operations Research and Management Science. Springer, 2008.
- [14] J. Mayer. *Stochastic Linear Programming Algorithms: A Comparison Based on a Model Management System*. Gordon and Breach Science Publishers, 1998.
- [15] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19:1574–1609, 2009.
- [16] 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.

- [17] A. Prékopa. *Stochastic Programming*. Kluwer Academic Publishers, Dordrecht, 1995.
- [18] A. Prékopa, B. Vizvári, and T. Badics. Programming under probabilistic constraint with discrete random variable. In F. Giannesi, T. Rapcsák, and S. Komlósi, editors, *New Trends in Mathematical Programming*, pages 235–255. Kluwer, Dordrecht, 1998.
- [19] H. Robbins and S. Monro. A stochastic approximation method. *Annals of Mathematical Statistics*, 22:400–407, 1951.
- [20] A. Ruszczyński. *Nonlinear Optimization*. Princeton University Press, 2006.
- [21] T. Szántai. *Numerical Evaluation of Probabilities Concerning Multidimensional Probability Distributions, Thesis*. Hungarian Academy of Sciences, Budapest, 1985.
- [22] 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.
- [23] 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.