

A stochastic approximation method for chance-constrained nonlinear programs

Rohit Kannan¹ and James Luedtke²

¹Wisconsin Institute for Discovery, University of Wisconsin-Madison, Madison, WI, USA.
E-mail: rohit.kannan@wisc.edu

²Department of Industrial & Systems Engineering and Wisconsin Institute for Discovery, University of Wisconsin-Madison, Madison, WI, USA. E-mail: jim.luedtke@wisc.edu

December 19, 2018

Abstract

We propose a stochastic approximation method for approximating the efficient frontier of chance-constrained nonlinear programs. Our approach is based on a bi-objective viewpoint of chance-constrained programs that seeks solutions on the efficient frontier of optimal objective value versus risk of constraints violation. In order to be able to apply a projected stochastic subgradient algorithm to solve our reformulation with the probabilistic objective, we adapt existing smoothing-based approaches for chance-constrained problems to derive a convergent sequence of smooth approximations of our reformulated problem. In contrast with exterior sampling-based approaches (such as sample average approximation) that approximate the original chance-constrained program with one having finite support, our proposal converges to local solutions of a smooth approximation of the original problem, thereby avoiding poor local solutions that may be an artefact of a fixed sample. Computational results on three test problems from the literature indicate that our proposal is consistently able to determine better approximations of the efficient frontier than existing approaches in reasonable computation times. We also present a bisection approach for solving chance-constrained programs with a prespecified risk level.

Key words: stochastic approximation, chance constraints, efficient frontier, stochastic subgradient

1 Introduction

We consider the solution of the following class of chance-constrained nonlinear programs (NLPs):

$$\begin{aligned} \nu^* &:= \min_{x \in X} f(x) \\ \text{s.t. } &\mathbb{P}\{g(x, \xi) \leq 0\} \geq 1 - \alpha, \end{aligned} \tag{CCP}$$

where $X \subset \mathbb{R}^n$ is a nonempty closed convex set, ξ is a random vector with probability distribution \mathcal{P} supported on $\Xi \subset \mathbb{R}^d$, $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a continuous quasiconvex function, $g : \mathbb{R}^n \times \mathbb{R}^d \rightarrow \mathbb{R}^m$ is continuously differentiable, the constant $\alpha \in (0, 1)$ is a user-defined acceptable level of constraints violation, and for each $x \in \mathbb{R}^n$, we write $\mathbb{P}\{g(x, \xi) \leq 0\}$ to denote the probability with which the random constraint $g(x, \xi) \leq 0$ is satisfied. Our contribution is to develop theory and a practical implementation of a stochastic approximation algorithm for approximating the efficient frontier of optimal objective value ν^* versus risk level α for (CCP). Most of our results only require mild assumptions on the data defining the above problem such as the ability to draw i.i.d. samples of ξ from \mathcal{P} (see Section 3). In particular, we will not impose restrictive assumptions such as Gaussianity of the distribution \mathcal{P} or make structural assumptions on the random constraint functions g . Our result regarding ‘convergence to stationary solutions’ of (CCP) additionally requires ξ to be a continuous random variable satisfying some distributional assumptions. We note that Problem (CCP) can model joint nonconvex chance constraints, nonconvex deterministic constraints (by incorporating them in g), and even some recourse structure (by defining g through the solution of auxiliary optimization problems, see Appendix C).

Chance-constrained programming was introduced as a modeling framework for optimization under uncertainty by Charnes et al. [14], and was soon after generalized to the joint chance-constrained case by Miller and Wagner [40] and to the nonlinear case by Prékopa [44]. Apart from a few known tractable cases, e.g., Prékopa [45] and Lagoa et al. [35], solving chance-constrained NLPs is in general hard since the feasible region is not guaranteed to be convex (a notable case is when the underlying deterministic problem is itself nonconvex) and even merely evaluating the probabilistic constraint involves multi-dimensional integration. Motivated by a diverse array of applications [10, 29, 37, 63], most existing numerical approaches for chance-constrained NLPs attempt to determine good-quality feasible solutions in reasonable computation times (see Section 2 for details).

The last decade has also seen a surge of interest in another subarea of stochastic optimization that encompasses optimization problems with nonconvex expected value objectives and simple convex constraints, primarily inspired by applications in machine learning. This has resulted in a flurry of advances, especially in the development and analysis of first-order algorithms, e.g., [8, 19, 22, 23, 26, 27, 42], that guarantee finding first-order stationary points of a broad class of such problems in a well-defined sense. Inspired by the success of these first-order methods, we explore the option of using them to approximate the efficient frontier of (CCP) in this work. We begin with the following hypothesis:

In many cases, decision makers are interested in generating the efficient frontier of optimal objective function value (ν^*) versus risk level (α) rather than simply solving (CCP) for a single prespecified risk level so that they can make a more informed decision [38, 47].

In this work, we investigate the potential of a stochastic subgradient algorithm [19, 27] for approximating the efficient frontier of (CCP). While our proposal is more naturally applicable for approximating the efficient frontier, we also present a bisection strategy for solving (CCP) for a fixed risk level α .

We begin by observing that the efficient frontier can be ‘recovered’ by solving the following stochastic optimization problem [47]:

$$\begin{aligned} \min_{x \in X} \quad & \mathbb{P}\{g(x, \xi) \leq 0\} \\ \text{s.t.} \quad & f(x) \leq \nu, \end{aligned} \tag{1}$$

where the above formulation determines the minimum probability of constraints violation given an upper bound ν on the objective function value. Problem (1) is conceivably easier to handle than (CCP) because the random variables only participate in the objective. Additionally, the above reformulation crucially enables the use of stochastic subgradient algorithms for its approximate solution as will be demonstrated shortly. Assuming that (CCP) is feasible for risk levels of interest, it is not hard to see that solving (1) to global optimality using different appropriate values of the bound ν will yield the same efficient frontier as solving (CCP) using different (corresponding) values of α [47].

Let $\mathbb{1} : \mathbb{R} \rightarrow \{0, 1\}$, defined by $\mathbb{1}[z] = 1$ if $z > 0$ and $\mathbb{1}[z] = 0$ if $z \leq 0$, denote the characteristic function of the set $(0, +\infty)$ (also called the lower semicontinuous Heaviside/step function), and $\max_j [\mathbb{1}[g_j(x, \xi)]]$ denote $\max_{j \in \{1, \dots, m\}} \mathbb{1}[g_j(x, \xi)]$. Note that the following problem is equivalent to (1):

$$\min_{x \in X_\nu} \mathbb{E}[\max_j [\mathbb{1}[g_j(x, \xi)]]],$$

where X_ν denotes the closed convex set $\{x \in X : f(x) \leq \nu\}$. The above formulation is almost in the form that stochastic gradient-type algorithms can handle [20], but poses at least a couple of challenges that (potentially) preclude the use of such algorithms for its solution: the step function $\mathbb{1}[\cdot]$ is discontinuous and the max function is nondifferentiable. We propose to solve a partially-smoothed approximation of the above formulation using the projected stochastic subgradient algorithm of Davis and Drusvyatskiy [19]. To enable this, we replace the step functions in the above formulation using smooth¹ approximations $\phi_j : \mathbb{R} \rightarrow \mathbb{R}$, $j = 1, \dots, m$, to obtain the following approximation to (1):

$$\min_{x \in X_\nu} \mathbb{E}[\max_j [\phi_j(g_j(x, \xi))]], \tag{APP}$$

¹We will refer to the approximations to the step function as ‘smooth’ throughout this work, although continuously differentiable approximations with Lipschitz continuous gradients will suffice for the purposes of our analysis. Note that in a departure from previous works [25, 53], we do not restrict our approximations of the step function to be conservative.

where $\max[\phi(g(x, \xi))]$ is shorthand for the composite function $\max[\phi_1(g_1(x, \xi)), \dots, \phi_m(g_m(x, \xi))]$. Since the objective function of (APP) is weakly convex under mild assumptions [22], using the projected stochastic subgradient algorithm of Davis and Drusvyatskiy [19] will guarantee convergence to a neighborhood of an approximately stationary solution \bar{x}_ν of (APP) for each well-chosen value of ν . The probability of constraints violation $\mathbb{P}\{g(\bar{x}_\nu, \xi) \not\leq 0\}$ can then be estimated at \bar{x}_ν to determine a point that is approximately on the efficient frontier. Because we rely on an algorithm that does not guarantee finding globally optimal solutions to (APP), our proposal may not yield the efficient frontier of (CCP), but may only result in a conservative (‘achievable’) approximation of it. Section 4 presents some options for the smooth approximations ϕ_j .

A major motivation for our developments is the fact that sample average approximations [39] of (CCP) may introduce spurious local optima as a byproduct of sampling, see Curtis et al. [18, Fig. 1]. This may unnecessarily complicate the process of obtaining a good approximate solution to (CCP), especially in cases when ξ is a ‘well-behaved’ continuous random variable that makes the probability function $p(x) := \mathbb{P}\{g(x, \xi) \not\leq 0\}$ ‘well behaved’ (even though it might be nonconvex). While smoothing approaches [25] help mitigate this issue, their practical deployment still requires integration within a sample average approximation (SAA) framework. By employing a stochastic subgradient algorithm, we do not restrict ourselves to a single fixed batch of samples of the random variables, and can therefore hope to converge to truly locally optimal solutions of (APP). Another motivating factor is that theoretical bounds [11, 39] on the number of samples required by SAA-based approaches to ensure feasibility for (CCP) are typically very conservative [31, 39, 43], and their use results in suboptimal solutions and increased solution times (this issue is alleviated through the use of ‘tuned SAA’ and carefully designed solution approaches, see Section 6). Our proposal, on the other hand, leverages well-known advantages of stochastic subgradient-type approaches, including a low per-iteration cost, low memory requirement, good scaling with number of random variables and the number of joint chance constraints, and reasonable scaling with the number of decision variables.

An alternative to approximating the solution of (1) using (APP) is to solve Problem (1) directly using zeroth-order stochastic optimization approaches, see Ghadimi et al. [27, Section 5] and Jin et al. [34], for instance. A potential disadvantage of such approaches is that they may require taking a relatively large sample (mini-batch) of ξ at each iteration to make reasonable progress, since they ignore the structure of the objective function of (1) altogether. As an aside, we note that our proposal can also be used for obtaining near-stationary solutions of classes of reliability maximization problems [33, 46, 50] since such problems can be cast directly in the form (1).

Notation. We use e to denote a vector of ones of suitable dimension, $\mathbb{1}[\cdot]$ to denote the l.s.c. step function, $X_\nu := \{x \in X : f(x) \leq \nu\}$, $p(x) := \mathbb{P}\{g(x, \xi) \not\leq 0\}$, $B_\delta(x) := \{z \in \mathbb{R}^n : \|z - x\| < \delta\}$, where $\|\cdot\|$ denotes the Euclidean norm, and ‘a.e.’ for the quantifier ‘almost everywhere’ with respect to the probability measure \mathbb{P} . Given functions $h : \mathbb{R}^n \times \mathbb{R}^d \rightarrow \mathbb{R}^m$ and $\phi_j : \mathbb{R} \rightarrow \mathbb{R}$, $j = 1, \dots, m$, we write $\max[h(x, \xi)]$ to denote $\max_{j \in \{1, \dots, m\}} h_j(x, \xi)$ and $\phi[h(x, \xi)]$ to denote the element-wise composition $(\phi_1(h_1(x, \xi)), \dots, \phi_m(h_m(x, \xi)))$. Given a set $Z \subset \mathbb{R}^n$, we let $\text{co}\{Z\}$ denote its convex hull, $I_Z(\cdot)$ denote its characteristic function, i.e., $I_Z(z) = +\infty$ if $z \notin Z$, and zero otherwise, $N_Z(z)$ to denote the normal cone of Z at a vector $z \in \mathbb{R}^n$ and $\text{proj}(z, Z)$ to denote its projection onto Z (we abuse notation to write $y = \text{proj}(z, Z)$ when Z is a nonempty closed convex set), and $|Z|$ to denote its cardinality when it is a finite set. Given a twice differentiable function $f : \mathbb{R} \rightarrow \mathbb{R}$, we write $df(x)$ and $d^2f(x)$ to denote its first and second derivatives at $x \in \mathbb{R}$. Given a locally Lipschitz continuous function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, we write $\partial f(x)$ to denote its Clarke generalized gradient at a point $x \in \mathbb{R}^n$ [16] and $\partial_B f(x)$ to denote its B -subdifferential (note: $\partial f(x) = \text{co}\{\partial_B f(x)\}$).

This article is organized as follows. Section 2 reviews related approaches for solving chance-constrained NLPs. Section 3 establishes consistency of the smoothing approach in the limit of its parameter values, and Section 4 lists some choices for smooth approximation of the step function. Section 5 outlines our proposal for approximating the efficient frontier of (CCP) and lists some practical considerations for its implementation. Section 6 summarizes implementation details and presents computational experiments on three test problems that illustrate the strength of our approach. We close with a summary of our contributions and some avenues for future work in Section 7. The appendices present auxiliary algorithms, proofs, and computational results, and provide technical and implementation-related details for problems with recourse structure.

2 Review of related work

Several algorithms have been proposed in the literature for ‘approximately solving’ chance-constrained programs. In this section, we restrict our attention to algorithms that attempt to generate good-quality *feasible* solutions to chance-constrained *nonlinear* programs.

The conceptually simplest (but remarkably versatile) approach, termed *scenario approximation*, takes a fixed sample $\{\xi^i\}_{i=1}^N$ of the random variables from the distribution \mathcal{P} and enforces the sampled constraints $g(x, \xi^i) \leq 0$, $\forall i \in \{1, \dots, N\}$, in lieu of the chance constraint $\mathbb{P}\{g(x, \xi) \leq 0\} \geq 1 - \alpha$. The advantage of this approach is that the scenario problem is a standard NLP that can be solved using off-the-shelf solvers. Calafiore, Campi, and Garatti [9, 11] present upper bounds on the sample size N required to ensure that the solution of the scenario problem is *feasible* for (CCP) with a high user-specified probability when the functions g are partly convex with respect to x for each fixed $\xi \in \Xi$. When the functions g do not possess the above convexity property but instead satisfy alternative Lipschitz continuity assumptions, Luedtke and Ahmed [39, Section 2.2.3] determine an upper bound on N that provides similar theoretical guarantees for a slightly perturbed version of the scenario problem. While scenario approximation has been a popular approach for obtaining *feasible* solutions to (CCP), its major downside is that it does not provide guarantees on solution quality without a significant increase in computational effort [39, 47].

There has been a recent surge of interest in iterative smoothing-based NLP approaches in which the hard chance constraint $p(x) \leq \alpha$ is replaced by a convergent sequence of ‘easier, smooth constraints’ for approximately solving (CCP). To illustrate these approaches, consider first the case when we only have an individual chance constraint (i.e., $m = 1$), and suppose $p(x) = \mathbb{E}[\mathbb{1}[g(x, \xi)]] \leq \alpha$ is approximated by $\mathbb{E}[\phi_k(g(x, \xi))] \leq \alpha$, where $\{\phi_k\}$ is a monotonically ‘convergent’ sequence of smooth approximations of the step function. If each element of $\{\phi_k\}$ is chosen to overestimate the step function, then solving the sequence of smooth approximations will furnish feasible solutions to (CCP) of improving quality. Since there are no known general approaches for solving optimization problems with nonconvex expectation constraints (the algorithm of Lan and Zhou [36] can be used to solve problems with convex expectation constraints), the smooth approximation is solved in practice using a scenario-based approach. Joint chance constraints can be accommodated either by replacing the max function in $\mathbb{E}[\max\{\phi_k(g(x, \xi))\}]$ with its own convergent sequence of smooth approximations [52], or by conservatively approximating (CCP) using a Bonferroni-type approximation [41]. We describe some iterative smoothing-based approaches for (CCP) below. As an aside, we mention that feasible solutions for (CCP) could also be obtained using conservative convex optimization-based approaches [15, 41, 48] when the constraint functions g are partly convex with respect to x and satisfy additional structural assumptions; however, their application may result in overly conservative solutions, e.g., see Hong et al. [32] and Cao and Zavala [12].

Hong et al. [32] propose a sequence of conservative nonsmooth difference-of-convex (DC) approximations of the step function and use it to solve joint chance-constrained convex programs. They establish convergence of any stationary/optimal solution of the DC approximation scheme to the set of stationary/optimal solutions of (CCP) in the limit under certain assumptions. Shan et al. [52, 53] build on the above work by developing a family of conservative smooth DC approximations of the step function, employ conservative log-sum-exp approximations of the max function to construct smooth approximations of (CCP), and provide similar theoretical guarantees as the above work under slightly weaker assumptions. Geletu et al. [25] propose a framework for conservative analytic approximations of the step function, provide similar theoretical guarantees as the above two works, and illustrate the applicability of their framework for individual chance-constrained NLPs using a class of sigmoid-like smoothing functions. They also propose smooth outer-approximations of (CCP) for generating lower bounds. Adam et al. [1, 3] develop a continuous relaxation of SAAs of (CCP), and propose to use smooth approximations of the step function that are borrowed from the literature on mathematical programs with complementarity constraints along with Benders’ decomposition [6] to determine stationary points. Finally, Cao and Zavala [12] propose a nonsmooth sigmoidal approximation of the step function, and use it to solve a smoothed SAA of individual chance-constrained NLPs. They propose to initialize their algorithm using the solution of a scenario approximation of a conservative convex optimization problem [48].

An alternative to smoothing the chance constraint is to solve a SAA problem directly using mixed-integer nonlinear programming techniques, which may be computationally cumbersome, especially when the constraint functions g are not partly convex with respect to x . Several tailored approaches [38, 57] have been proposed for solving SAAs of chance-constrained convex programs in a bid to reduce

computational burden. Curtis et al. [18] attempt to directly solve a SAA of (CCP) using nonlinear programming techniques. They develop an exact penalty function for the SAA, and propose a trust region algorithm that solves quadratic programs with linear cardinality constraints to converge to stationary points of (CCP).

Another important body of work by Henrion et al. [55, 56] establishes (sub)gradient formulae for the probability function p when the function g possesses special structure and \mathcal{P} is a Gaussian or Gaussian-like distribution. These approaches employ internal sampling to numerically evaluate (sub)gradients, and can be used within a nonlinear programming framework for computing stationary point solutions to (CCP). When Problem (CCP) possesses certain special structures, nonlinear programming approaches based on ‘ p -efficient points’ [21, 58] can also be employed for its solution.

The recent independent work of Adam and Branda [2] also considers a stochastic approximation method for chance-constrained NLPs. They consider the case when the random vector ξ has a discrete distribution with finite (but potentially large) support, rewrite the chance constraint in (CCP) using a quantile-based constraint, develop a penalty-based approach to transfer this constraint to the objective, and employ a mini-batch projected stochastic subgradient method to determine an approximately stationary solution to (CCP) for a given risk level. Because their manuscript does not contain any theory justifying the applicability and convergence of their proposal, it currently only offers a heuristic approach.

The major distinguishing feature of our proposal from the above works is that it does not solve a SAA of (CCP) that relies on a fixed sample of the random variables, but instead attempts to solve (1) using a stochastic approximation algorithm that uses implicit sampling. Consequently, our proposed algorithm typically has a low per-iteration computational cost that is well-suited for large-scale instances, and can hope to converge to local optimal solutions of (1). This is demonstrated by the numerical experiments in Section 6 which indicate that our approach can find good-quality approximations of the efficient frontier in reasonable solution times. Our proposal is unfortunately not entirely immune to drawbacks; some of them are listed at the end of Section 5 after introducing our proposed algorithm formally.

3 Consistency of the smoothing approach

In this section, we establish sufficient conditions under which solving a sequence of smooth approximations (APP) of the stochastic program (1) yields a point on the efficient frontier of (CCP). While the techniques used in this section are not new (cf. Hong et al. [32], Shan et al. [52, 53], and Geletu et al. [25]), our analysis is necessitated by the following reasons: i. our approximations (APP) involve a convergent sequence of approximations of the objective function of (1) rather than a convergent sequence of approximations of the feasible region of (CCP), ii. we handle joint chance-constrained NLPs directly without ‘reducing’ them to the case of individual chance-constrained programs, and iii. our computational results rely on a slight generalization of existing frameworks for smooth approximation of the step function. We will make the following assumptions on Problem (1) for every $\nu \in \mathbb{R}$ of interest, *as needed*, to ensure convergence of optimal objective values and solutions of our approximation scheme.

Assumption 1. The convex set $X_\nu := \{x \in X : f(x) \leq \nu\}$ is nonempty and compact.

Assumption 2. For each $x \in X_\nu$, we have $\mathbb{P}\{\max [g(x, \xi)] = 0\} = 0$.

Assumption 3. The following conditions on the functions g_j hold for each $\xi \in \Xi$ and $j \in \{1, \dots, m\}$:

- (3a) The function $g_j(\cdot, \xi)$ is Lipschitz continuous on X_ν with a nonnegative measurable Lipschitz constant $L_{g,j}(\xi)$ satisfying $\mathbb{E}[L_{g,j}^2(\xi)] < +\infty$.
- (3b) The gradient function $\nabla g_j(\cdot, \xi)$ is Lipschitz continuous on X_ν with a measurable Lipschitz constant $L'_{g,j}(\xi) \in \mathbb{R}_+$ satisfying $\mathbb{E}[L'_{g,j}(\xi)] < +\infty$.
- (3c) There exist positive constants $\sigma_{g,j}$ satisfying $\mathbb{E}[\|\nabla_x g_j(\bar{x}, \xi)\|^2] \leq \sigma_{g,j}^2, \forall \bar{x} \in X_\nu$.

Assumption 1 is used to ensure that Problem (1) is well defined (see Lemma 1). Assumption 2 guarantees that the family of approximations of the function p considered converge pointwise to it on X_ν in the limit of their parameter values (see Proposition 1). This assumption, when enforced, pretty much restricts ξ to be a continuous random variable (see Lemma 1 for a consequence of this assumption).

Assumption 3 serves two purposes: it ensures that the approximations of p possess important regularity properties (see Lemma 2), and it enables the use of the projected stochastic subgradient algorithm for solving the sequence of approximations (APP). Note that Assumption 2 is implied by Assumption (6b), which will be introduced later, and that Assumption (3a) implies $\mathbb{E}[L_{g,j}(\xi)] < +\infty$.

Lemma 1. The probability function $p : \mathbb{R}^n \rightarrow [0, 1]$ is lower semicontinuous. Furthermore, under Assumption 2, p is continuous on X_ν .

Proof. Follows from Theorem 10.1.1 of Prékopa [45]. \square

A natural approach to approximating the solution of Problem (1) (and one that we will justify in Section 5) is to construct a sequence of approximations (APP) based on an associated sequence of smooth approximations that converge to the step function (cf. Section 2). In what follows, we use $\{\phi_k\}$ to denote a sequence of smooth approximations of the step function, where ϕ_k corresponds to a vector of approximating functions $(\phi_{k,1}, \dots, \phi_{k,m})$ for the m constraints defined by the function g . For ease of exposition, we make the following (mild) blanket assumptions on each element of the sequence of smoothing functions $\{\phi_k\}$ throughout this work. Section 4 lists some examples that satisfy these assumptions.

Assumption 4. The following conditions hold for each $k \in \mathbb{N}$ and $j \in \{1, \dots, m\}$:

- (4a) The functions $\phi_{k,j} : \mathbb{R} \rightarrow \mathbb{R}$ are continuously differentiable.
- (4b) Each function $\phi_{k,j}$ is nondecreasing, i.e., $y \geq z \implies \phi_{k,j}(y) \geq \phi_{k,j}(z)$.
- (4c) The functions $\phi_{k,j}$ are equibounded, i.e., there exists a universal constant $M_\phi > 0$ (independent of indices j and k) such that $|\phi_{k,j}(y)| \leq M_\phi, \forall y \in \mathbb{R}$.
- (4d) Each approximation $\phi_{k,j}(y)$ converges pointwise to the step function $\mathbf{1}[y]$ except possibly at $y = 0$, i.e., $\lim_{k \rightarrow \infty} \phi_{k,j}(y) = \mathbf{1}[y], \forall y \in \mathbb{R} \setminus \{0\}$.

We say that ‘the *strong form* of Assumption 4’ holds if Assumption (4d) is replaced with the stronger condition of pointwise convergence *everywhere*, i.e., $\lim_{k \rightarrow \infty} \phi_{k,j}(y) = \mathbf{1}[y], \forall y \in \mathbb{R}$. Note that a sequence of conservative smooth approximations of the step function (which overestimate the step function everywhere) cannot satisfy the strong form of Assumption 4, whereas sequences of underestimating smooth approximations of the step function may satisfy it (see Section 4). In the rest of this paper, we will use $\hat{p}_k(x)$ to denote the approximation $\mathbb{E}[\max\{\phi_k(g(x, \xi))\}]$. Note that $\hat{p}_k : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuous under Assumption 4, see Shapiro et al. [54, Theorem 7.43]. The following result establishes sufficient conditions under which $\hat{p}_k(x) \rightarrow p(x)$ pointwise on X_ν , which is a minimum requirement for the consistency of the smoothing approach in general.

Proposition 1. Suppose Assumptions 2 and 4 (or only the strong form of Assumption 4) hold. Then $\lim_{k \rightarrow \infty} \hat{p}_k(x) = p(x), \forall x \in X_\nu$.

Proof. See Appendix B.1. \square

The above result establishes the equivalence of (APP) and (1) when the smoothing parameters approach their limiting values, and indicates that approximating Problem (1) using a good approximation of the step function in Problem (APP) might not be a terrible idea. The next two results show that a global solution of (1) can be obtained by solving a sequence of approximations (APP) to global optimality. To achieve this, we rely on the related concept of epi-convergence of sequences of extended real-valued functions (see Rockafellar and Wets [49, Chapter 7] for an introduction).

Proposition 2. Suppose Assumptions 2 and 4 hold. Then, the sequence of functions $\{\hat{p}_k(\cdot) + I_{X_\nu}(\cdot)\}_k$ epi-converges to $p(\cdot) + I_{X_\nu}(\cdot)$ for each $\nu \in \mathbb{R}$.

Proof. See Appendix B.2. \square

A consequence of the above proposition is the following key result (cf. Hong et al. [32, Theorem 2], Shan et al. [53, Theorem 4.1], Geletu et al. [25, Corollary 3.7], and Cao and Zavala [12, Theorem 5]), which establishes convergence of the optimal solutions and objective values of the sequence of approximating problems (APP) to those of the true problem (1).

Theorem 1. Suppose Assumptions 1, 2, and 4 hold. Then

$$\lim_{k \rightarrow \infty} \min_{x \in X_\nu} \hat{p}_k(x) = \min_{x \in X_\nu} p(x) \quad \text{and} \quad \limsup_{k \rightarrow \infty} \arg \min_{x \in X_\nu} \hat{p}_k(x) \subset \arg \min_{x \in X_\nu} p(x).$$

Proof. Follows from Proposition 2 and Theorem 7.33 of Rockafellar and Wets [49]. \square

The above result, while important, has at least a couple of practical limitations. Firstly, it is not applicable to situations where ξ is a discrete random variable since it relies crucially on Assumption 2. Secondly, it only establishes that any accumulation point of a sequence of *global* minimizers of (APP) is a *global* minimizer of (1). Since (APP) involves minimizing a nonsmooth nonconvex expected-value function, this is not a practically useful guarantee. The next result circumvents the first limitation when the smoothing function is chosen judiciously (cf. the outer-approximations of Geletu et al. [25] and Section 4). Proposition 4 shows that strict local minimizers of (1) can be approximated using sequences of local minimizers of (APP). While this doesn't fully address the second limitation, it provides some hope that strict local minimizers of (1) may be approximated by solving a sequence of approximations (APP) to local optimality. Theorem 2 establishes that accumulation points of sequences of stationary solutions to the approximations (APP) yield stationary solutions to (1) under additional assumptions.

Proposition 3. Suppose the approximations $\{\phi_k\}$ satisfy the strong form of Assumption 4. Then the sequence of functions $\{\hat{p}_k(\cdot) + I_{X_\nu}(\cdot)\}_k$ epi-converges to $p(\cdot) + I_{X_\nu}(\cdot)$. Moreover, the conclusions of Theorem 1 hold if we further make Assumption 1.

Proof. Follows by noting that Assumption 2 is no longer required for the proof of Proposition 2 when the strong form of Assumption 4 is made. \square

In fact, the proof of Proposition 3 becomes straightforward if we impose the additional (mild) condition (assumed by Geletu et al. [25]) that the sequence of smooth approximations $\{\phi_k\}$ is monotone nondecreasing, i.e., $\phi_{k+1,j}(y) \geq \phi_{k,j}(y)$, $\forall y \in \mathbb{R}, k \in \mathbb{N}$, and $j \in \{1, \dots, m\}$. Under this additional assumption, the conclusions of Proposition 3 follow from Proposition 7.4(d) and Theorem 7.33 of Rockafellar and Wets [49]. The next result, in the spirit of Proposition 3.9 of Geletu et al. [25], shows that strict local minimizers of (1) can be approximated using local minimizers of (APP).

Proposition 4. Suppose the assumptions of Theorem 1 (or the assumptions of Proposition 3) hold. If x^* is a strict local minimizer of (1), then there exists a sequence of local minimizers $\{\hat{x}_k\}$ of $\min_{x \in X_\nu} \hat{p}_k(x)$ with $\hat{x}_k \rightarrow x^*$.

Proof. See Appendix B.3. \square

Note that we make do with the above result because we are unable to establish the more desirable statement that a convergent sequence of local minimizers of approximations (APP) converges to a local minimizer of (1) without additional assumptions (cf. Theorem 2). The next few results work towards establishing conditions under which a convergent sequence of stationary solutions of the sequence of approximating problems (APP) converges to a stationary point of (1). We will make the following additional assumptions on each element of the sequence of smoothing functions $\{\phi_k\}$ for this purpose.

Assumption 5. The following conditions hold for each $k \in \mathbb{N}$ and $j \in \{1, \dots, m\}$:

- (5a) The derivative mapping $d\phi_{k,j}(\cdot)$ is bounded by $M'_{\phi,k,j}$ on \mathbb{R} , i.e., $|d\phi_{k,j}(y)| \leq M'_{\phi,k,j}$, $\forall y \in \mathbb{R}$.
- (5b) The derivative mapping $d\phi_{k,j}(\cdot)$ is Lipschitz continuous on \mathbb{R} with Lipschitz constant $L'_{\phi,k,j} \in \mathbb{R}_+$.

The above assumptions are mild since we let the constants $M'_{\phi,k,j}$ and $L'_{\phi,k,j}$ depend on the sequence index k . Note that in light of Assumption (4a), Assumption (5a) is equivalent to the assumption that $\phi_{k,j}(\cdot)$ is Lipschitz continuous on \mathbb{R} with Lipschitz constant $M'_{\phi,k,j}$. We make the following assumptions on the cumulative distribution function of a 'scaled version' of the constraint functions g and on the sequence $\{\phi_{k,j}\}$ of smoothing functions to establish Proposition 5 and Theorem 2.

Assumption 6. The following conditions on the constraint functions g , the distribution of ξ , and the sequences $\{\phi_{k,j}\}$ of smoothing functions hold for each $k \in \mathbb{N}$:

- (6a) There exist positive constants β_1, \dots, β_m such that the approximations $\hat{p}_k(x) := \mathbb{E}[\max[\phi_k(g(x, \xi))]]$ can be rewritten as $\hat{p}_k(x) = \mathbb{E}[\bar{\phi}_k(\max[\bar{g}(x, \xi)])]$, where $\bar{g} : \mathbb{R}^n \times \mathbb{R}^d \rightarrow \mathbb{R}^m$ is defined by $\bar{g}_j(x, \xi) := \beta_j g_j(x, \xi)$, $\forall (x, \xi) \in \mathbb{R}^n \times \mathbb{R}^d$ and $j \in \{1, \dots, m\}$, and $\bar{\phi}_k$ is a scalar smoothing function satisfying Assumptions 4 and 5.
- (6b) Let $F : \mathbb{R}^n \times \mathbb{R} \rightarrow [0, 1]$ be the cumulative distribution function of $\max[\bar{g}(x, \xi)]$, i.e., $F(x, \eta) := \mathbb{P}\{\max[\bar{g}(x, \xi)] \leq \eta\}$. There exists a constant $\theta > 0$ such that distribution function F is continuously differentiable on $\bar{X} \times (-\theta, \theta)$, where $\bar{X} \supset X_\nu$ is an open subset of \mathbb{R}^n . Furthermore, for each $\eta \in \mathbb{R}$, $F(\cdot, \eta)$ is Lipschitz continuous on X_ν with a measurable Lipschitz constant $L_F(\eta) \in \mathbb{R}_+$ that is also Lebesgue integrable, i.e., $\int L_F(\eta) d\eta < +\infty$.
- (6c) There exists a sequence of positive constants $\{\varepsilon_k\} \downarrow 0$ such that

$$\lim_{k \rightarrow \infty} \int_{|\eta| \geq \varepsilon_k} L_F(\eta) d\bar{\phi}_k(\eta) d\eta = 0 \quad \text{and} \quad \lim_{k \rightarrow \infty} \bar{\phi}_k(\varepsilon_k) - \bar{\phi}_k(-\varepsilon_k) = 1.$$

Some comments on the above assumption are in order. Assumption (6a) is a mild assumption on the choice of smoothing functions ϕ_k that essentially assumes that the approximations \hat{p}_k can be formulated in one of two equivalent ways: either using a tailored smoothing function $\phi_{k,j}$ for each constraint function g_j , or by using the same smoothing function ϕ_k on a fixed rescaling \bar{g}_j of the constraint functions g . Note that all three examples of the sequence of smoothing functions $\{\phi_{k,j}\}$ in Section 4 satisfy this assumption in a typical implementation. The Lipschitz continuity assumption in Assumption (6b) is mild (and is similar to Assumption 4 of Shan et al. [53]). Assumption (6b), which also assumes local continuous differentiability of the distribution function F , is quite strong (this assumption is similar to Assumption 4 of Hong et al. [32]). A consequence of this assumption is that the function $p(x) \equiv 1 - F(x, 0)$ is continuously differentiable on \bar{X} . Finally, note that Assumption (6c) is mild, see Section 4 for examples that satisfy it. When made along with Assumption 4, the first part of this assumption ensures that the derivative mapping $d\bar{\phi}_k(\cdot)$ approaches the ‘Dirac delta function’ sufficiently rapidly.

The following result ensures that the Clarke gradient of the approximation \hat{p}_k is well defined.

Lemma 2. Suppose Assumptions 3, 4, and 5 hold. Then for each $k \in \mathbb{N}$:

1. $\phi_{k,j}(g_j(\cdot, \xi))$ is Lipschitz continuous on X_ν with Lipschitz constant $M'_{\phi,k,j} L_{g,j}(\xi)$ for each $\xi \in \Xi$ and $j \in \{1, \dots, m\}$.
2. \hat{p}_k is Lipschitz continuous on X_ν with Lipschitz constant $\mathbb{E} \left[\max_{j \in \{1, \dots, m\}} \left[M'_{\phi,k,j} L_{g,j}(\xi) \right] \right]$.

Proof. See Appendix B.4. □

The next result characterizes the Clarke generalized gradient of the approximating objectives \hat{p}_k .

Proposition 5. Suppose Assumptions 3, 4, and 5 hold. Then the Clarke generalized gradient of \hat{p}_k can be expressed as

$$\partial \hat{p}_k(x) = \mathbb{E}[\text{co}\{\Gamma_k(x, \xi)\}],$$

where $\Gamma_k : \mathbb{R}^n \times \mathbb{R}^d \rightrightarrows \mathbb{R}^n$ is defined as $\Gamma_k(x, \xi) := \{\nabla_x \phi_{k,l}(g_l(x, \xi)) : l \in \mathcal{A}(x, \xi)\}$, $\mathcal{A}(x, \xi)$ denotes the set of active indices $l \in \{1, \dots, m\}$ at which $\max[\phi_k(g(x, \xi))] = \phi_{k,l}(g_l(x, \xi))$, and the expectation above is to be interpreted in the sense of Definition 7.39 of Shapiro et al. [54].

Proof. See Appendix B.5. □

The next result, similar to Lemma 4.2 of Shan et al. [53], helps characterize the accumulation points of stationary solutions to the sequence of approximations (APP).

Proposition 6. Suppose Assumptions 3, 4, 5, and 6 hold. Then

$$\limsup_{\substack{x \rightarrow \bar{x} \\ k \rightarrow \infty}} \partial \hat{p}_k(x) + N_{X_\nu}(x) \subset \nabla p(\bar{x}) + N_{X_\nu}(\bar{x}).$$

Proof. See Appendix B.6. □

The following key result is an immediate consequence of the above proposition (cf. Shan et al. [53, Theorem 4.2]), and establishes convergence of Clarke stationary solutions of the sequence of approximating problems (APP) to those of the true problem (1).

Theorem 2. Suppose Assumptions 1, 3, 4, 5, and 6 hold. Let $\{x_k\}$ be a sequence of stationary solutions to the sequence of approximating problems $\min_{x \in X_\nu} \hat{p}_k(x)$. Then every accumulation point of $\{x_k\}$ is an stationary solution to $\min_{x \in X_\nu} p(x)$.

Proof. From the corollary to Proposition 2.4.3 of Clarke [16], we have $0 \in \partial \hat{p}_k(x_k) + N_{X_\nu}(x_k)$ by virtue of the stationarity of x_k for $\min_{x \in X_\nu} \hat{p}_k(x)$. The stated result then follows from Proposition 6 by noting that $0 \in \nabla p(\bar{x}) + N_{X_\nu}(\bar{x})$ for any accumulation point \bar{x} of $\{x_k\}$ (see Theorem 5.37 of [49]). \square

4 Examples of smooth approximations

We present a few examples of sequences of smooth approximations of the step function that satisfy Assumptions 4, 5, and 6. Throughout this section, we let $\{\tau_{k,j}\}$, $j \in \{1, \dots, m\}$, denote sequences of *positive* reals with $\lim_k \tau_{k,j} = 0$, $\forall j \in \{1, \dots, m\}$.

Example 1. This example is from Example 5.1 of Shan et al. [53]. The sequence of approximations

$$\phi_{k,j}(y) = \begin{cases} 0 & \text{if } y < -\tau_{k,j}, \\ 1 - 2 \left(\frac{y}{\tau_{k,j}} \right)^3 - 3 \left(\frac{y}{\tau_{k,j}} \right)^2 & \text{if } -\tau_{k,j} \leq y \leq 0, \\ 1 & \text{if } y > 0 \end{cases}$$

of the step function satisfy Assumptions 4 and 5. Assumption (6a) is easily verified under Equation (2), and Assumption (6c) trivially holds with $\varepsilon_k = \tau_k$.

Example 2. This is based on the above example and the outer-approximation framework of Geletu et al. [25]. The sequence of approximations

$$\phi_{k,j}(y) = \begin{cases} 0 & \text{if } y < 0, \\ 1 - 2 \left(\frac{y - \tau_{k,j}}{\tau_{k,j}} \right)^3 - 3 \left(\frac{y - \tau_{k,j}}{\tau_{k,j}} \right)^2 & \text{if } 0 \leq y \leq \tau_{k,j}, \\ 1 & \text{if } y > \tau_{k,j} \end{cases}$$

of the step function satisfy the strong form of Assumption 4 and Assumption 5. Assumption (6a) is easily verified under Equation (2), and Assumption (6c) trivially holds with $\varepsilon_k = \tau_k$.

The next example is the sequence of smooth approximations $\{\phi_k\}$ that we adopt in this work.

Example 3. The sequence of approximations

$$\phi_{k,j}(y) = \frac{1}{1 + \exp\left(-\frac{y}{\tau_{k,j}}\right)}$$

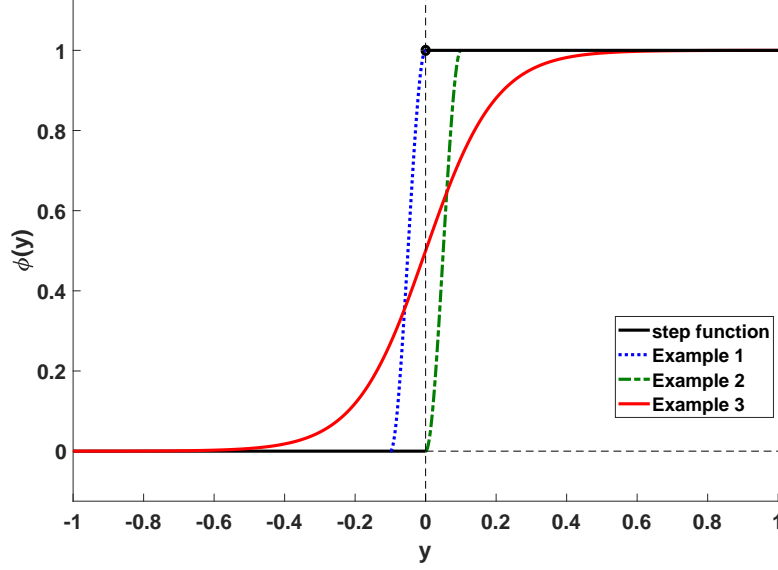
of the step function satisfy Assumptions 4 and 5 (see Proposition 7). Assumption (6a) is easily verified under Equation (2), and Assumption (6c) could be satisfied with $\varepsilon_k = \tau_k^c$ for some constant $c \in (0, 1)$ (depending on the distribution function F).

Figure 1 illustrates the above smoothing functions with the common smoothing parameter $\tau = 0.1$. We refer the reader to the works of Shan et al. [53], Geletu et al. [25], and Cao and Zavala [12] for other examples of ‘smooth’ approximations that can be accommodated within our framework. While the results established in Section 3 provide a great deal of flexibility in choosing the sequence of smoothing parameters $\{\tau_k\}$, we restrict our implementation to the geometric setting

$$\tau_{k+1,j} = \tau_c \tau_{k,j}, \quad \forall k \in \mathbb{N}, j \in \{1, \dots, m\} \quad (2)$$

for some constant factor $\tau_c \in (0, 1)$. The next result estimates some important constants related to our smooth approximation in Example 3.

Figure 1: Illustration of the examples of smoothing functions in Section 4.



Proposition 7. The sequence of approximations $\{\phi_k\}$ in Example 3 satisfies Assumption 5 with constants $M'_{\phi,k,j} \leq 0.25\tau_{k,j}^{-1}$ and $L'_{\phi,k,j} \leq 0.1\tau_{k,j}^{-2}$, $\forall k \in \mathbb{N}$, $j \in \{1, \dots, m\}$.

Proof. See Appendix B.7. □

5 Proposed algorithm

Our proposal for approximating the efficient frontier of (CCP) involves solving a sequence of approximations (APP) constructed using the sequence of smooth approximations of the step function introduced in Example 3. We use the projected stochastic subgradient algorithm of Davis and Drusvyatskiy [19] to obtain an approximately stationary solution to (APP) for each value of the objective bound ν and each element of the sequence of smoothing parameters $\{\tau_k\}$ (note that we suppress the constraint index j in our notation from hereon unless necessary). Section 5.1 outlines a conceptual algorithm for approximating the efficient frontier of (CCP) using the projected stochastic subgradient algorithm, and Sections 5.2 and 5.3 outline our proposal for estimating some parameters of the algorithm in Section 5.1 that are important for a good practical implementation. We close this section with a discussion of our proposed approach in Section 5.4.

5.1 Outline of the algorithm

Algorithm 1 outlines our proposal for approximately discretizing the efficient frontier of (CCP). This algorithm takes as its inputs an initial guess $\hat{x}^0 \in X$, an initial objective bound $\bar{\nu}^0$ of interest, an objective spacing $\tilde{\nu}$ for discretizing the efficient frontier, and a lower bound α_{low} on risk levels of interest that is used as a termination criterion. Section 5.2 prescribes ways to determine a good initial guess \hat{x}^0 and initial bound $\bar{\nu}^0$, whereas Section 6.1 lists our setting for $\tilde{\nu}$, α_{low} , and other algorithmic parameters. The preprocessing step uses Algorithm 3 to determine a suitable initial sequence of smoothing parameters $\{\tau_k\}$ based on the problem data, and uses this choice of $\{\tau_k\}$ to determine an initial sequence of step lengths for the stochastic subgradient algorithm (see Section 5.2). It is worth mentioning that good choices of both of the above parameters are critical for our proposal to work well. The optimization phase of Algorithm 1 uses Algorithm 2 to solve the sequence of approximations (APP) for different choices of the objective bound $\bar{\nu}^i$ (note that Algorithm 2 rescales its input smoothing parameters $\{\bar{\tau}_k\}$ at

the initial guess \bar{x}^i for each objective bound $\bar{\nu}^i$). Finally, Algorithm 7 in the appendix adapts Algorithm 1 to solve (CCP) approximately for a given risk level $\hat{\alpha}$.

Algorithm 1 Approximating the efficient frontier of (CCP)

- 1: **Input:** initial point $\hat{x}^0 \in X$, initial objective bound $\bar{\nu}^0$, objective spacing $\tilde{\nu}$, and lower bound on risk levels of interest $\alpha_{low} \in (0, 1)$.
 - 2: **Set algorithmic parameters:** mini-batch size M , maximum number of iterations N_{max} , minimum and maximum number of ‘runs’, R_{min} and R_{max} , parameters for determining and updating step length, termination criteria, fixed sample $\{\bar{\xi}^l\}_{l=1}^{N_{MC}}$ from \mathcal{P} for estimating risk levels of candidate solutions, and sequence of smoothing parameters $\{\bar{\tau}_k\}_{k=1}^K$ for some $K \in \mathbb{N}$.
 - 3: **Output:** set of pairs $\{(\bar{\nu}^i, \bar{\alpha}^i)\}$ of objective values and risk levels that can be used to approximate the efficient frontier and associated solutions $\{\hat{x}^i\}$.
 - 4: **Preprocessing:** determine smoothing parameters $\{\tau_k\}_{k=1}^K$ scaled at $\text{proj}(\hat{x}^0, X_{\bar{\nu}^0})$ using Algorithm 3, and an initial sequence of step lengths $\{\bar{\gamma}_k\}_{k=1}^K$ for the corresponding sequence of approximating problems (APP) using Algorithm 4.
 - 5: **Optimization Phase:**
 - 6: Initialize iteration count $i = 0$.
 - 7: **repeat**
 - 8: Update iteration count $i \leftarrow i + 1$ and set objective bound $\bar{\nu}^i = \bar{\nu}^0 - (i - 1)\tilde{\nu}$.
 - 9: Obtain $(\bar{\alpha}^i, \hat{x}^i)$ by solving the sequence of approximations (APP) using Algorithm 2 with scaling parameters $\{\tau_k\}_{k=1}^K$, the above algorithmic parameter settings, and $\bar{x}^i := \text{proj}(\hat{x}^{i-1}, X_{\bar{\nu}^i})$ as the initial guess.
 - 10: **until** $\bar{\alpha}^i \leq \alpha_{low}$
-

Algorithm 2 solves a sequence of approximations (APP) for a given objective bound ν using the projected stochastic subgradient algorithm of Davis and Drusvyatskiy [19] and an adaptation of the two-phase randomized stochastic projected gradient method of Ghadimi et al. [27]. This algorithm begins by rescaling the smoothing parameters $\{\bar{\tau}_k\}$ at the initial guess \bar{x}^0 using Algorithm 3 to ensure that the sequence of approximations (APP) are well-scaled at \bar{x}^0 . The optimization phase then solves each element of the sequence of approximations (APP) using two loops: the inner loop employs a mini-batch version² of the algorithm of Davis and Drusvyatskiy [19] to solve (APP) for a given smoothing parameter τ_k , and the outer loop assesses the progress of the inner loop across multiple ‘runs/replicates’ (the initialization strategy for each run is based on Algorithm 2-RSPG-V of Ghadimi et al. [27]). Finally, the initial guess for the next approximating problem (APP) is set to be a solution corresponding to the smallest estimate of the risk level determined thus far. This initialization step is important for the next approximating problem in the sequence to be well-scaled at its initial guess - this is why we do not solve a ‘single tight approximating problem’ (APP) that can be hard to initialize (cf. Figure 1), but instead solve a sequence of approximations to approximate the solution of (1).

In the remainder of this section, we verify that the assumptions of Davis and Drusvyatskiy [19] hold to justify the use of their projected stochastic subgradient method for solving the sequence of approximations (APP). First, we verify that the objective function \hat{p}_k of (APP) with smoothing parameter τ_k is a weakly convex function on X_ν (see Definition 4.1 of Drusvyatskiy and Paquette [22]). Next, we verify that Assumption (A3) of Davis and Drusvyatskiy [19] holds.

Proposition 8. Suppose Assumptions 3, 4 and 5 hold. Then $\hat{p}_k(\cdot)$ is \bar{L}_k -weakly convex on X_ν , where \bar{L}_k is the Lipschitz constant of the Jacobian $\mathbb{E}[\nabla \phi_k(g(\cdot, \xi))]$ on X_ν .

Proof. Follows from Lemma 4.2 of Drusvyatskiy and Paquette [22]. □

The following result will prove useful for bounding \hat{p}_k ’s weak convexity parameter in terms of known constants.

Proposition 9. Suppose Assumptions 3, 4 and 5 hold. Then

²Although Davis and Drusvyatskiy [19] establish that mini-batching is not necessary for the convergence of the projected stochastic subgradient algorithm, a small mini-batch can greatly enhance the performance of the algorithm in practice.

Algorithm 2 Generating a point approximately on the efficient frontier

```

1: Input: objective bound  $\nu$ , initial guess  $\bar{x}^0 \in X_\nu$ , mini-batch size  $M$ , maximum number of iterations
    $N_{max}$ , minimum and maximum number of ‘runs’,  $R_{min}$  and  $R_{max} \geq R_{min}$ , initial step lengths
    $\{\bar{\gamma}_k\}_{k=1}^K$ , parameters for updating step length, replicate termination criteria, sample  $\{\bar{\xi}^l\}_{l=1}^{N_{MC}}$  from
    $\mathcal{P}$  for estimating risk levels, and sequence of smoothing parameters  $\{\bar{\tau}_k\}_{k=1}^K$ .
2: Output: Smallest estimate of the risk level  $\hat{\alpha}$  and corresponding solution  $\hat{x} \in X_\nu$ .
3: Preprocessing: determine smoothing parameters  $\{\tau_k\}_{k=1}^K$  scaled at  $\bar{x}^0$  using Algorithm 3.
4: Initialize: Best found solution  $\hat{x} = \bar{x}^0$  and estimate of its risk level  $\hat{\alpha}$  using the sample  $\{\bar{\xi}^l\}_{l=1}^{N_{MC}}$ .
5: Optimization Phase:
6: for approximation  $k = 1$  to  $K$  do
7:   Initialize run count  $i = 0$  and step length  $\gamma_k = \bar{\gamma}_k$ .
8:   repeat
9:     Update  $i \leftarrow i + 1$ , and initialize the first iterate  $x^1 := \bar{x}^{i-1}$ .
10:    Choose number of iterations  $N$  uniformly at random from  $\{1, \dots, N_{max}\}$ .
11:    for iteration  $l = 1$  to  $N - 1$  do
12:      Draw an i.i.d. sample  $\{\xi^{l,q}\}_{q=1}^M$  of  $\xi$  from  $\mathcal{P}$ .
13:      Estimate an element  $G(x^l) \in \mathbb{R}^n$  of the subdifferential of the objective of (APP) with
        smoothing parameters  $\tau_k$  at  $x^l$  using the mini-batch  $\{\xi^{l,q}\}_{q=1}^M$ .
14:      Update  $x^{l+1} := \text{proj}(x^l - \gamma_k G(x^l), X_\nu)$ .
15:    end for
16:    Set  $\bar{x}^i = x^N$ , and estimate its risk level  $\bar{\alpha}^i$  using the sample  $\{\bar{\xi}^l\}_{l=1}^{N_{MC}}$ .
17:    Update step length  $\gamma_k$  using Algorithm 5, and the incumbents  $(\hat{\alpha}, \hat{x})$  with  $(\bar{\alpha}^i, \bar{x}^i)$  if  $\hat{\alpha} > \bar{\alpha}^i$ .
18:  until run termination criteria are satisfied (see Algorithm 5)
19:  Update initial guess  $\bar{x}^0 = \hat{x}$  for the next approximation.
20: end for

```

1. For any $\xi \in \Xi$, the Lipschitz constant $L_{k,j}(\xi)$ of $\nabla \phi_{k,j}(g_j(\cdot, \xi))$ on X_ν satisfies

$$L_{k,j}(\xi) \leq M'_{\phi,k,j} L'_{g,j}(\xi) + L'_{\phi,k,j} L_{g,j}^2(\xi).$$

2. The Lipschitz constant \bar{L}_k of the Jacobian $\nabla \mathbb{E}[\phi_k(g(\cdot, \xi))]$ on X_ν satisfies

$$\bar{L}_k \leq \mathbb{E} \left[\left(\sum_{j=1}^m L_{k,j}^2(\xi) \right)^{\frac{1}{2}} \right].$$

Proof. See Appendix B.8. □

We derive alternative results regarding the weak convexity parameter of the approximation \hat{p}_k when (CCP) is used to model recourse formulations in Appendix C. We now establish that Assumption (A3) of Davis and Drusvyatskiy [19] holds.

Proposition 10. Suppose Assumptions 3, 4, and 5 hold. Let $G(x, \xi)$ be a stochastic B -subdifferential element of $\hat{p}_k(\cdot)$ at $x \in X_\nu$, i.e., $\mathbb{E}[G(x, \xi)] \in \partial_B \mathbb{E}[\max[\phi_k(g(x, \xi))]]$. Then

$$\mathbb{E}[\|G(x, \xi)\|^2] \leq \max_{j \in \{1, \dots, m\}} \left(M'_{\phi,k,j} \right)^2 \sigma_{g,j}^2.$$

Proof. Follows from Assumptions (3c) and 5 and the fact that $G(x, \xi) = d\phi_{k,j}(g_j(x, \xi)) \nabla_x g_j(x, \xi)$ for some active constraint index $j \in \{1, \dots, m\}$ such that $\max[\phi_k(g(x, \xi))] = \phi_{k,j}(g_j(x, \xi))$. □

5.2 Estimating the parameters of Algorithm 1

This section outlines our proposal for estimating some key parameters of Algorithm 1. We present pseudocode for determining suitable smoothing parameters and step lengths, and for deciding when to terminate Algorithm 2 with a point on our approximation of the efficient frontier.

Algorithm 3 rescales the sequence of smoothing parameters $\{\bar{\tau}_k\}$ based on the values assumed by the constraints g at a reference point \bar{x} and a Monte Carlo sample of the random vector ξ . The purpose of this rescaling step is to ensure that the initial approximation (APP) with smoothing parameter τ_1 is well-scaled at the initial guess \bar{x} . If the initial guess \bar{x} for the above approximation is ‘near’ a stationary solution, then we can hope that the approximation remains well-scaled over a region of interest. Note that the form of the scaling factor β_j in Algorithm 3 is influenced by our choice of the smoothing functions in Example 3. Algorithm 4 uses the step length rule in Davis and Drusvyatskiy [19, Page 6] with the trivial bound ‘ $R = 1$ ’ and sample estimates of the weak convexity parameter of \hat{p}_k and the parameter σ_k^2 related to its stochastic subgradient. Since stochastic approximation algorithms are infamous for their sensitivity to the choice of step lengths (see Section 2.1 of Nemirovski et al. [42], for instance), Algorithm 5 prescribes heuristics for updating the step length based on the progress of Algorithm 2 over multiple runs. These rules increase the step length if insufficient progress has been made, and decrease the step length if things have gone significantly downhill in the last few runs. Algorithm 5 also suggests the following heuristic for terminating the replicate loop within Algorithm 2: terminate either if the upper limit on the number of runs is hit, or if insufficient progress has been made over the last few runs.

Algorithm 3 Scaling the smoothing parameters

- 1: **Input:** reference point $\bar{x} \in X$, and sequence of smoothing parameters $\{\bar{\tau}_k\}_{k=1}^K$.
 - 2: **Set algorithmic parameters:** number of samples N_{scale} , scaling tolerance $s_{tol} > 0$, and scaling factor $\omega > 0$.
 - 3: **Output:** smoothing parameters $\{\tau_k\}_{k=1}^K$ scaled at \bar{x} .
 - 4: Draw a sample $\{\xi^i\}_{i=1}^{N_{scale}}$ of ξ from \mathcal{P} .
 - 5: For each $j \in \{1, \dots, m\}$ and k , set $\tau_{k,j} = \beta_j \bar{\tau}_{k,j}$, where $\beta_j := \omega \max \{ \text{median} \{ |g_j(\bar{x}, \xi^i)| \}, s_{tol} \}$.
-

Algorithm 4 Determining initial step length

- 1: **Input:** objective bound ν , reference point $\bar{x} \in X_\nu$, mini-batch size M , maximum number of iterations N_{max} , minimum number of runs R_{min} , and smoothing parameters $\{\tau_k\}_{k=1}^K$.
 - 2: **Set algorithmic parameters:** number of samples for estimating weak convexity parameter N_{wc} , sampling radius $r > 0$, number of samples for estimating ‘variability’ of stochastic gradients N_{var} , and batch size for computing expectations N_{batch} .
 - 3: **Output:** initial step lengths $\{\bar{\gamma}_k\}_{k=1}^K$.
 - 4: Estimate the weak convexity parameter ρ_k of the objective of (APP) with smoothing parameter τ_k using i. N_{wc} samples of pairs of points in $X_\nu \cap B_r(\bar{x})$, and ii. mini-batches of ξ for estimating stochastic gradients (see Section 5.3).
 - 5: Estimate $\sigma_k^2 := \max_{x \in X_\nu} \mathbb{E} [\|G(x, \xi)\|^2]$ using i. N_{var} sample points $x \in X_\nu \cap B_r(\bar{x})$, and ii. stochastic B -subdifferential elements $G(x, \xi)$ of the subdifferential of the objective of (APP) with smoothing parameter τ_k and mini-batches of ξ (see Section 5.3).
 - 6: Set $\bar{\gamma}_k := \frac{1}{\sqrt{\rho_k \sigma_k^2 (N_{max} + 1) R_{min}}}$.
-

5.3 Other practical considerations

We list some other practical considerations for implementing Algorithm 1 below.

Setting an initial guess: The choice of the initial point \hat{x}^0 for Algorithm 1 is important especially because Algorithm 2 does not guarantee finding global solutions to (APP). Additionally, a poor choice of the initial bound $\bar{\nu}^0$ can lead to excessive effort expended on uninteresting regions of the efficient frontier. We propose to initialize \hat{x}^0 and $\bar{\nu}^0$ by solving a set of tuned scenario approximation problems. For instance, Algorithm 6 in the appendix can be solved with $M = 1$, $R = 1$, and a tuned sample size N_1 (tuned such that the output $(\bar{\nu}^{1,1}, \bar{\alpha}^{1,1}, \bar{x}^{1,1})$ of this algorithm corresponds to an interesting region of the efficient frontier) to yield the initial guess $\hat{x}^0 = \bar{x}^{1,1}$ and $\bar{\nu}^0 = \bar{\nu}^{1,1}$. Note that it may be worthwhile to solve the scenario approximation problems to global optimality to obtain a good initialization $(\hat{x}^0, \bar{\nu}^0)$,

Algorithm 5 Updating step length and checking termination

```

1: Input: minimum and maximum number of ‘runs’,  $R_{min}$  and  $R_{max}$ , step length  $\gamma_k$ , and estimated
   risk levels at the end of each run thus far.
2: Set algorithmic parameters: step length checking frequency  $N_{check}$ , replicate termination check-
   ing parameter  $N_{term}$ , relative increase in risk level  $\delta_1 > 0$  beyond which we terminate, relative
   increase in risk level  $\delta_2 > \delta_1$  beyond which we decrease the step length by factor  $\gamma_{decr} > 1$ , relative
   decrease in risk level  $\delta_3 > 0$  beyond which we increase the step length by factor  $\gamma_{incr} > 1$ .
3: Updating step length:
4: if run number is divisible by  $N_{check}$  then
5:   Determine maximum relative ‘decrease’  $\hat{\delta} \leq 1$  in the estimated risk level over the past  $N_{check}$ 
     runs, relative to the smallest known estimate of the risk level  $N_{check}$  runs ago.
6:   if  $\hat{\delta} \in (-\delta_1, \delta_3)$  then
7:     Increase step length by factor  $\gamma_{incr}$ . ▷ Insufficient decrease in estimated risk level
8:   else if  $\hat{\delta} \leq -\delta_2$  then
9:     Decrease step length by factor  $\gamma_{decr}$ . ▷ Unacceptable increase in estimated risk level
10:  end if
11: end if
12: Termination check:
13: if number of runs equals  $R_{max}$  then
14:   Terminate.
15: else if number of runs is at least  $R_{min}$  then
16:   Determine maximum relative ‘decrease’  $\hat{\delta} \leq 1$  in the estimated risk level over the past
      $N_{term}$  runs. If  $\hat{\delta} < -\delta_1$ , terminate due to insufficient decrease in risk level.
17: end if

```

e.g., when the number of scenarios N_1 is not too large (in which case they may be solved using off-the-shelf solvers/tailored decomposition techniques in acceptable computation times). Case study 3 in Section 6 provides an example where globally optimizing the scenario approximation problems to obtain an initial guess yields a significantly better approximation of the efficient frontier.

Computing projections: Algorithm 1 implicitly assumes that it is easy to project onto the sets X_ν for each $\nu \in \mathbb{R}$ of interest since it may involve several such steps while executing Algorithm 2. It may be worthwhile to implement tailored projection subroutines [13], especially for sets X_ν with special structures [17]. If the set X_ν is a polyhedron, then projection onto X_ν may be carried out by solving a quadratic program.

Computing stochastic subgradients: Given a point $\bar{x} \in X_\nu$ and a mini-batch $\{\xi^l\}_{l=1}^M$ of the random vector ξ , a stochastic element $G(\bar{x}, \xi)$ of the B -subdifferential of the objective of (APP) with smoothing parameter τ_k at \bar{x} can be obtained as

$$G(\bar{x}, \xi) = \frac{1}{M} \sum_{l=1}^M d\phi_{k,j(l)}(g_{j(l)}(\bar{x}, \xi^l)) \nabla_x g_{j(l)}(\bar{x}, \xi^l),$$

where for each $l \in \{1, \dots, M\}$, $j(l) \in \{1, \dots, m\}$ denotes a constraint index satisfying $\max [\phi_k(g(\bar{x}, \xi^l))] = \phi_{k,j(l)}(g_{j(l)}(\bar{x}, \xi^l))$. It is worth mentioning that sparse computation of stochastic subgradients can speed Algorithm 1 up significantly.

Estimating risk levels: Given a candidate solution $\bar{x} \in X_\nu$ and a sample $\{\bar{\xi}^l\}_{l=1}^{N_{MC}}$ of the random vector ξ , a stochastic upper bound on the risk level $p(\bar{x})$ can be obtained as

$$\bar{\alpha} = \max_{\alpha \in [0,1]} \left\{ \alpha : \sum_{i=0}^{N_{viol}} \binom{N_{MC}}{i} \alpha^i (1-\alpha)^{N_{MC}-i} = \delta \right\},$$

where N_{viol} is the cardinality of the set $\{l : g(\bar{x}, \bar{\xi}^l) \not\leq 0\}$ and $1 - \delta$ is the required confidence level, see Nemirovski and Shapiro [41, Section 4]. Since checking the satisfaction of N_{MC} constraints at the

end of each run in Algorithm 2 may be time consuming (especially for problems with recourse structure), we use a smaller sample size (determined based on the risk lower bound α_{low}) to estimate risk levels during the course of Algorithm 2 and only use all N_{MC} samples of ξ to estimate the risk level $\{\bar{\alpha}^i\}$ of the final solutions $\{\hat{x}^i\}$ in Algorithm 1. If our discretization of the efficient frontier obtained from Algorithm 1 consists of N_{EF} points each of whose risk levels is estimated using a confidence level of $1 - \delta$, then we may conclude that our approximation of the efficient frontier is ‘achievable’ with a confidence level of at least $1 - \delta N_{EF}$ using Bonferroni’s inequality.

Estimating weak convexity parameters: Proposition 8 shows that the Lipschitz constant \bar{L}_k of the Jacobian $\mathbb{E}[\nabla\phi_k(g(\cdot, \xi))]$ on X_ν provides a conservative estimate of the weak convexity parameter ρ_k of \hat{p}_k on X_ν . Therefore, we use an estimate of the Lipschitz constant \bar{L}_k as an estimate of ρ_k . To avoid overly conservative estimates, we restrict the estimation of \bar{L}_k to a neighborhood of the reference point \bar{x} in Algorithm 4 to get at the local Lipschitz constant of the Jacobian $\mathbb{E}[\nabla\phi_k(g(\cdot, \xi))]$.

Estimating σ_k^2 : For each point x sampled from X_ν , we compute multiple realizations of the mini-batch stochastic subdifferential element $G(x, \xi)$ outlined above to estimate $\mathbb{E}[\|G(x, \xi)\|^2]$. Once again, we restrict the estimation of σ_k^2 to a neighborhood of the reference point \bar{x} in Algorithm 4 to get at the local ‘variability’ of the stochastic subgradients.

Estimating initial step lengths: Algorithm 4 estimates initial step lengths for the sequence of approximations (APP) by estimating the parameters $\{\rho_k\}$ and $\{\sigma_k^2\}$, which in turn involve estimating subgradients and the Lipschitz constants of $\{\hat{p}_k\}$. Since the numerical conditioning of the approximating problems (APP) deteriorates rapidly as the smoothing parameters τ_k approach zero (see Figure 1), obtaining good estimates of these constants via sampling becomes tricky when k is ‘large’. To circumvent this difficulty, we only estimate the initial step length $\bar{\gamma}_1$ for the initial approximation (APP) with smoothing parameter τ_1 by sampling, and propose the conservative initialization $\bar{\gamma}_k := \left(\frac{\tau_k}{\tau_1}\right)^2 \bar{\gamma}_1$ for $k > 1$ (see Propositions 7, 8, 9, and 10 for a justification).

5.4 Discussion of the proposed approach

We close this section with a discussion of shortcomings and potential adaptations of our proposal.

Generalizing our approach to the case of multiple sets of joint chance constraints is nontrivial. For ease of exposition, suppose we have the $|\mathcal{J}|$ sets of joint chance constraints $\mathbb{P}\{g_j(x, \xi) \leq 0, \forall j \in J\} \geq 1 - \alpha_J$, $\forall J \in \mathcal{J}$, instead of the single joint chance constraint $\mathbb{P}\{g_j(x, \xi) \leq 0, \forall j \in J\} \geq 1 - \alpha$. If the decision maker wishes to impose distinct risk levels α_J for each set of joint chance constraints $J \in \mathcal{J}$, the natural analogue of our biobjective viewpoint would necessitate constructing the efficient frontier of a multiobjective optimization problem, which is considerably harder. If we assume that the risk levels are the same across the different sets of joint chance constraints (which may be a reasonable assumption in practice), we can write down the following analogue of Problem (APP) that is a harder-to-solve stochastic compositional minimization problem [61]:

$$\min_{x \in X_\nu} \max_{J \in \mathcal{J}} \left\{ \mathbb{E} \left[\max_{j \in J} [\phi_j(g_j(x, \xi))] \right] \right\}.$$

One potential avenue for reducing effort spent on projections in Algorithm 1 is to use a random constraint projection technique [60], although the associated theory will have to be developed for our nonsmooth nonconvex setting. Our proposal does not handle deterministic nonconvex constraints effectively. Although such constraints can theoretically be incorporated as part of the chance constraint functions g , this may not be a practically useful option, particularly when we have deterministic nonconvex equality constraints. A possibly practical option for incorporating deterministic nonconvex constraints in a more natural fashion is through carefully designed penalty-based stochastic approximation methods [62]. Note that Algorithm 1 does not accommodate discrete decisions as well.

Finally, while the theoretical results in Section 3 establish conditions under which solutions of the approximations (APP) converge to a solution of the stochastic program (1) in the limit of the smoothing parameters, Algorithm 1 only considered a finite sequence of approximating problems with smoothing

parameters $\{\tau_k\}_{k=1}^K$. This is because the numerical conditioning of the approximations (APP) deteriorates rapidly as the smoothing parameters approach their limiting value. Unfortunately, this seems to be a fundamental shortcoming of smoothing-based approaches in general since approximating the step function using a well-conditioned sequence of smooth approximations is impossible (cf. the discussion in Sections 4 and 5 of Cao and Zavala [12]).

6 Computational results

We present implementation details and results of our computational experiments in this section. We use the abbreviation ‘EF’ for the efficient frontier throughout this section.

6.1 Implementation details

The following parameter settings are used for testing our stochastic approximation method:

- Algorithm 1: Obtain initial guess \hat{x}^0 and initial objective bound \bar{v}^0 using Algorithm 6 in Appendix A.1 with $M = 1$, $R = 1$, and $N_1 = 10$. Set $\tilde{v} = 0.005|\bar{v}^0|$ and $\alpha_{low} = 10^{-4}$ unless otherwise specified. In addition, set $M = 20$, $N_{max} = 1000$, $R_{min} = 10$, $R_{max} = 50$, $K = 3$, and $\{\bar{\tau}_k\} = \{(0.1)^{k-1}\}$ (i.e., $\tau_c = 0.1$ in Equation (2)).
- Algorithm 3: $N_{scale} = 10^4$, $s_{tol} = 10^{-6}$, and $\omega = 1$.
- Algorithm 4: $N_{wc} = N_{var} = 200$, $N_{batch} = 20$, and $r = 0.1\|\bar{x}\|$.
- Algorithm 5: $N_{check} = 3$, $N_{term} = 5$, $\delta_1 = 10^{-4}$, $\delta_2 = 10^{-2}$, $\delta_3 = 10^{-1}$, and $\gamma_{incr} = \gamma_{decr} = 10$.
- Projecting onto X_ν : Case study 1: using the algorithm of Condat [17]; Case studies 2 and 3 and Case study 4 in the appendix: by solving a quadratic program.
- Estimating risk level $p(x)$: Case study 1: analytical solution; Case studies 2 and 3: Monte Carlo estimate using $N_{MC} = 10^5$ and reliability level $\delta = 10^{-6}$; Case study 4 in the appendix: numerical estimation of the exact risk level based on Ruben [51].
- Case studies 2 and 3: $\alpha_{low} = 5 \times 10^{-4}$ (since it is memory intensive to store a large Monte Carlo sample for Case study 2, and it is time intensive to estimate the risk level for a large Monte Carlo sample for Case study 3).

We compare the results of our proposed approach with a tuned application of scenario approximation that enforces a predetermined number of random constraints, solves the scenario problem to local/global optimality, and determines an a posteriori overestimate of the risk level at the solution of the scenario problem using an independent Monte Carlo sample (see Section 5.3) to estimate a point on the EF of (CCP). Algorithm 6 in Appendix A presents pseudocode for approximating the EF of (CCP) by solving a bunch of scenario problems using an iterative (cutting-plane) approach.

Our codes are written in Julia 0.6.2 [7], use Gurobi 7.5.2 [30] to solve linear, quadratic, and second-order cone programs, use IPOPT 3.12.8 [59] to solve nonlinear scenario approximation problems (with MUMPS [4] as the linear solver), and use SCIP 6.0.0 [28] to solve nonlinear scenario approximation problems to global optimality (if necessary). The above solvers were accessed through the JuMP 0.18.2 modeling interface [24]. All computational tests³ were conducted on a Surface Book 2 laptop running Windows 10 Pro with a 1.90 GHz four core Intel i7 CPU, 16 GB of RAM.

6.2 Numerical experiments

We tested our approach on three test cases from the literature. We present basic details of the test instances below. The reader is referred to <https://github.com/rohitkannan/SA-for-CCP> for Julia code and data of the test instances. For each test case, we compared the EFs generated by the stochastic approximation method against the solutions obtained using the tuned scenario approximation method and, if available, the analytical EF. Because the output of our proposed approach is random, we present enclosures of the EF generated by our proposal over ten different replicates in Appendix D for each case

³The scenario approximation problems solved to global optimality using SCIP in Case study 3 were run on a different (standard) laptop running Ubuntu 16.04 due to interfacing issues on the Windows platform.

study. The reader can verify that the output of the proposed method does not vary significantly across the replicates.

Case study 1. This portfolio optimization instance is based on Example 2.3.6 in Ben-Tal et al. [5], and includes a single individual linear chance constraint.

$$\begin{aligned} \max_{t, x \in \Delta_N} \quad & t \\ \text{s.t.} \quad & \mathbb{P}\{\xi^\top x \geq t\} \geq 1 - \alpha, \end{aligned}$$

where $\Delta_N := \{y \in \mathbb{R}_+^N : \sum_i y_i = 1\}$ is the standard simplex, x_i denotes the fraction of investment in stock $i \in \{1, \dots, N\}$, $\xi \sim \mathcal{P} := \mathcal{N}(\mu, \Sigma)$ is a random vector of returns with joint normal probability distribution, $\mu_i = 1.05 + 0.3 \frac{N-i}{N-1}$ and $\sigma_i = \frac{1}{3} \left(0.05 + 0.6 \frac{N-i}{N-1}\right)$, $i = 1, \dots, N$, and $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_N^2)$. We consider the instance with number of stocks $N = 1000$.

While our assumption that the returns ξ are normally distributed is not great, it allows us to benchmark our approach against an analytical solution for the true EF when the risk level $\alpha \leq 0.5$ (e.g., see Prékopa [45, Theorem 10.4.1]). Figure 2 compares a typical EF obtained using our approach against the analytical EF and the solutions generated by the tuned scenario approximation algorithm. Our proposal is able to find a very good approximation of the true EF, whereas the scenario approximation method only finds a poorer approximation. Our proposed approach took 267 seconds on average (and a maximum of 287 seconds) to approximate the EF using 26 points, whereas the tuned scenario approximation method took a total of 6900 seconds to generate its 1000 points in Figure 2.

Since existing smoothing-based approaches provide the most relevant comparison, we compare our results with one such approach from the literature. We choose the smoothing-based approach of Cao and Zavala [12] for comparison because: i. they report encouraging computational results in their work relative to other such approaches, and ii. they have made their implementation available. Table 1 summarizes typical results of the sigmoidal smoothing approach of Cao and Zavala [12] when applied to the above instance with a specified risk level of $\alpha = 0.01$, with a varying number of scenarios, and with different settings for the scaling factor γ of the sigmoidal approximation method (see Algorithm **SigVar-Alg** of Cao and Zavala [12]). Appendix A.3 lists details of our implementation of this method. The second, third, and fourth columns of Table 1 present the overall solution time in seconds (or a failure status returned by IPOPT), and the best objective value and risk level returned by the method over two replicates (we note that there was significant variability in solution times over the replicates; we report the results corresponding to the smaller solution times). Figure 2 plots the solutions returned by this method. The relatively poor performance of Algorithm **SigVar-Alg** on this example is not surprising; even the instance of the sigmoidal approximation problem with only a hundred scenarios (which is small for $\alpha = 0.01$) has more than a thousand variables and a hundred thousand nonzero entries in the (dense) Jacobian. Therefore, tailored approaches have to be developed for solving these problems efficiently. Since our implementation of the proposal of Cao and Zavala [12] failed to perform well on this instance even for a single risk level, we do not compare against their approach for the rest of the test instances.

Case study 2. This norm optimization instance is based on Section 5.1.2 of Hong et al. [32], and includes a joint convex chance constraint.

$$\begin{aligned} \min_{x \in \mathbb{R}_+^n} \quad & - \sum_i x_i \\ \text{s.t.} \quad & \mathbb{P}\left\{\sum_i \xi_{ij}^2 x_i^2 \leq U^2, \quad j = 1, \dots, m\right\} \geq 1 - \alpha, \end{aligned}$$

where ξ_{ij} are dependent normal random variables with mean $\frac{j}{d}$ and variance 1, and $\text{cov}(\xi_{ij}, \xi_{i'j}) = 0.5$ if $i \neq i'$, $\text{cov}(\xi_{ij}, \xi_{i'j'}) = 0$ if $j \neq j'$. We consider the instance with number of variables $n = 100$, number of constraints $m = 100$, and bound $U = 100$. Figure 3 compares a typical EF obtained using our approach against the solutions generated by the tuned scenario approximation algorithm. Once again, our proposed approach is able to find a significantly better approximation of the EF than the scenario approximation method, although it is unclear how our proposal fares compared to the true EF since it

Figure 2: Comparison of the efficient frontiers for Case study 1.

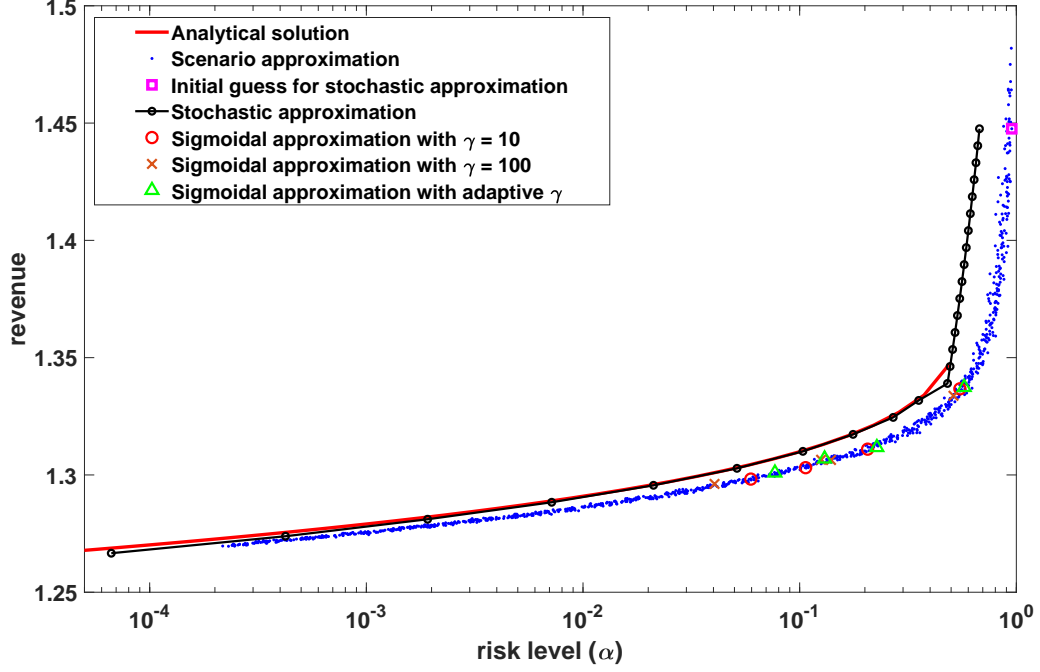
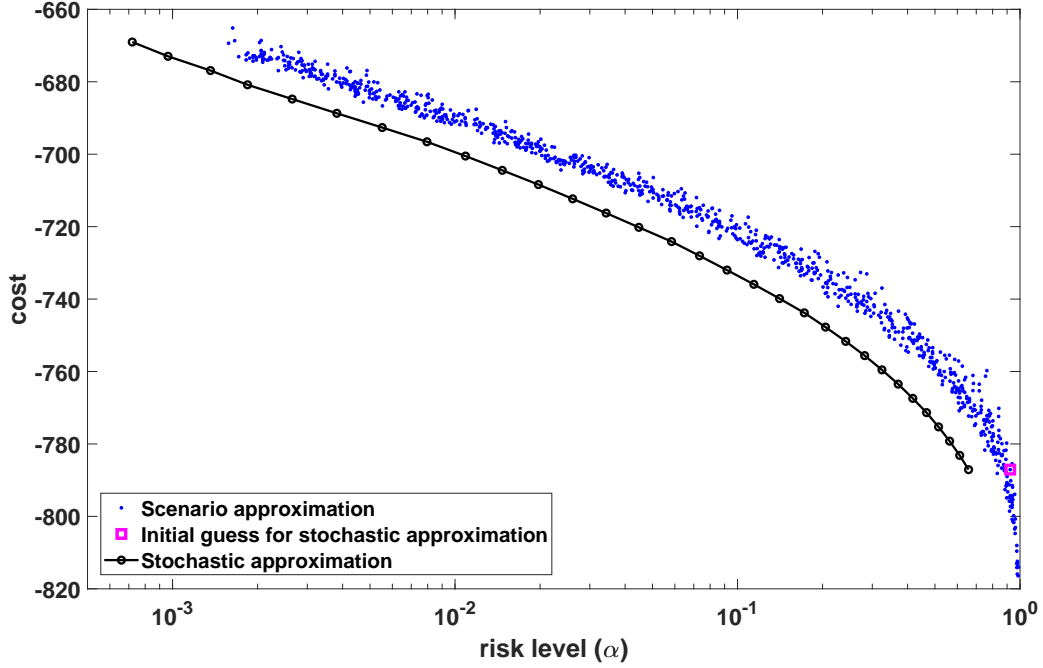


Table 1: Performance of the sigmoidal smoothing approach of Cao and Zavala [12] on Case study 1. The entry for each column of γ indicates the time in seconds (or IPOPT status)/best objective value/best risk level returned by the sigmoidal approximation method over two replicates. The entries **infeas** and **tle** denote infeasible and time limit exceeded statuses returned by IPOPT.

Num. scenarios	$\gamma = 10$			$\gamma = 100$			γ adaptive		
	Time	Obj	Risk	Time	Obj	Risk	Time	Obj	Risk
100	34	1.337	0.548	23	1.334	0.513	29	1.337	0.575
500	508	1.311	0.206	106	1.306	0.139	59	1.312	0.227
1000	606	1.303	0.106	1203	1.306	0.125	1011	1.307	0.130
2000	1946	1.298	0.059	296	1.296	0.040	2787	1.301	0.077
5000	4050	1.226	2.4×10^{-8}	6289	1.292	0.025	infeas		
10000	2708	1.199	5.1×10^{-11}	tle			tle		

Figure 3: Comparison of the efficient frontiers for Case study 2.



is unknown in this case. Our proposal took 5571 seconds on average (and a maximum of 5754 seconds) to approximate the EF using 31 points, whereas it took tuned scenario approximation a total of 25536 seconds to generate its 1000 points in Figure 3. We note that more than 70% of the reported times for our method is spent in generating random numbers because the random variable ξ is high-dimensional and the covariance matrix of the random vector $\xi_{\cdot j}$ is full rank. A practical instance might have a covariance matrix rank that is at least a factor of ten smaller, which would reduce our overall computation times roughly by a factor of three (cf. the smaller computation times reported for the similar Case study 4 in the appendix). Appendix D benchmarks our approach against the true EF when the random variables ξ_{ij} are assumed to be i.i.d. in which case the analytical solution is known (see Section 5.1.1 of Hong et al. [32]).

Case study 3. This probabilistic resource planning instance is based on Section 3 of Luedtke [38], and is modified to include a nonconvex recourse constraint.

$$\begin{aligned} \min_{x \in \mathbb{R}_+^n} \quad & c^T x \\ \text{s.t.} \quad & \mathbb{P}\{x \in R(\lambda, \rho)\} \geq 1 - \alpha, \end{aligned}$$

where x_i denotes the quantity of resource i , c_i denotes the unit cost of resource i ,

$$R(\lambda, \rho) = \left\{ x \in \mathbb{R}_+^n : \exists y \in \mathbb{R}_+^{n n_c} \text{ s.t. } \sum_{j=1}^{n_c} y_{ij} \leq \rho_i x_i^2, \forall i \in \{1, \dots, n\}, \sum_{i=1}^n \mu_{ij} y_{ij} \geq \lambda_j, \forall j \in \{1, \dots, n_c\} \right\},$$

n_c denotes the number of customer types, y_{ij} denote the amount of resource i allocated to customer type j , $\rho_i \in (0, 1]$ is a random variable that denotes the yield of resource i , $\lambda_j \geq 0$ is a random variable that denotes the demand of customer type j , and $\mu_{ij} \geq 0$ is a deterministic scalar that denotes the service rate of resource i for customer type j (we only let μ_{ij} be deterministic to enable an efficient implementation in JuMP). Note that the nonlinear term $\rho_i x_i^2$ in the definition of $R(\lambda, \rho)$ is a modification of the corresponding linear term $\rho_i x_i$ in Luedtke [38]. This change could be interpreted as a reformulation

of the instance in Luedtke [38] with concave objective costs (due to economies of scale). We consider the instance with number of resources $n = 20$ and number of customer types $n_c = 30$. Details of how the parameters of the model are set (including details of the random variables) can be found in the electronic companion to Luedtke [38].

Figure 4 compares a typical EF obtained using our approach against the solutions generated by the tuned scenario approximation algorithm. The blue dots in the top part of Figure 4 correspond to the 1000 points obtained using the scenario approximation method when IPOPT is used to solve the scenario approximation problems. We mention that the vertical axis has been re-scaled for readability; in reality, IPOPT returns solutions to the scenario approximation problem that are further away from the EF (with objective values up to 110 for the risk levels of interest). The 1000 red dots in the bottom part of Figure 4 correspond to the points obtained by solving the scenario approximation problems to global optimality using SCIP. The top black curve (with circles) corresponds to the EF obtained using the stochastic approximation method when it is initialized using the IPOPT solution of a scenario approximation problem. When a better initial point obtained by solving a scenario approximation problem using SCIP is used, the stochastic approximation method generates the bottom green curve (with squares) as its approximation of the EF.

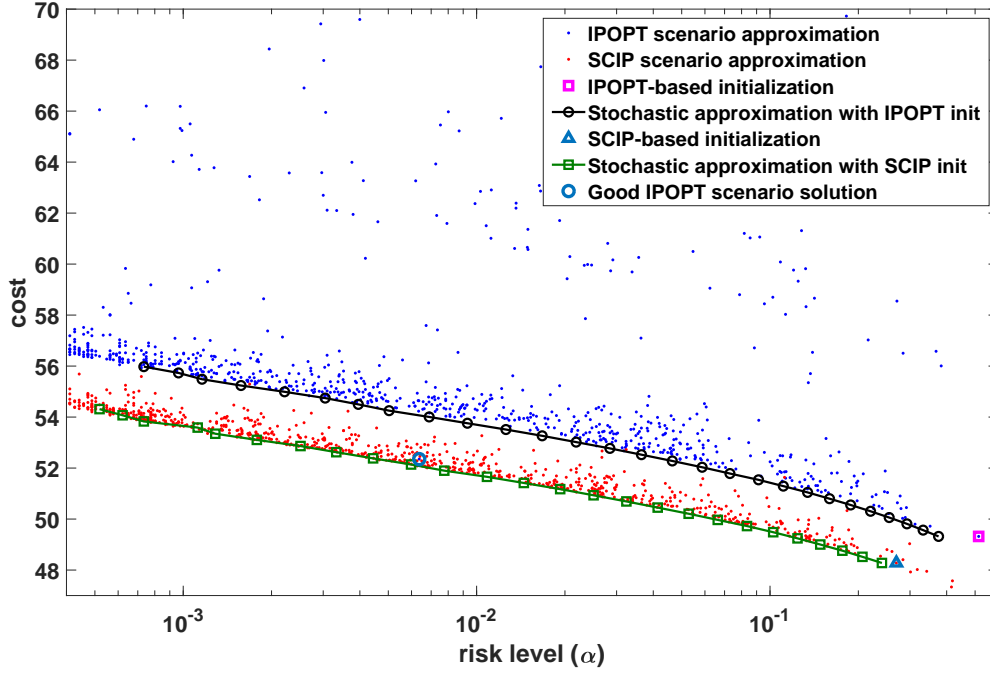
Several remarks are in order. The scenario approximation solutions generated using IPOPT are very scattered possibly because IPOPT gets stuck at suboptimal local minima. However, the best solutions generated using IPOPT provide a comparable approximation of the EF as the stochastic approximation method that is denoted by the black curve with circles (which appears to avoid the poor local minima encountered by IPOPT). We note that IPOPT finds a good local solution at one of the 1000 scenario approximation runs (indicated by the blue circle), which prompted us to use a global solver for solving the scenario approximation problems to initialize our algorithm and benchmark our solutions. Solving the scenario approximation problems to global optimality using SCIP yields a much better approximation of the EF than simply using IPOPT. Additionally, when the proposed approach is initialized using the SCIP solution of a single scenario approximation problem (which took less than 60 seconds to compute), it generates a significantly better approximation of the EF that performs comparably to the aforementioned approximation of the EF generated using SCIP. We mention that initializing the solution of IPOPT with the above initial guess from SCIP did not yield a better approximation of the EF than before - this may be because providing a good initialization to an interior point solver such as IPOPT is nontrivial. Our approach took 6978 seconds on average (and a maximum of 7356 seconds) to generate the green curve (with squares) approximation of the EF using 26 points, whereas it took the blue scenario approximations solved using IPOPT a total of 110291 seconds to generate its 1000 points in Figure 4.

7 Conclusion and future work

We proposed a stochastic approximation algorithm for estimating the efficient frontier of chance-constrained NLPs. Our proposal involves solving a sequence of partially smoothened stochastic optimization problems to local optimality using a projected stochastic subgradient algorithm. We established that every limit point of the sequence of stationary/global solutions of the above sequence of approximations yields a stationary/global solution of the original chance-constrained program with an appropriate risk level. A potential advantage of our proposal is that it can find truly stationary solutions of the chance-constrained NLP unlike scenario-based approaches that may get stuck at spurious local optima generated by sampling. Our computational experiments demonstrated that our proposed approach is consistently able to determine good approximations of the efficient frontier in reasonable computation times.

Extensions of our proposal that can handle multiple sets of joint chance constraints, for instance via the solution of the stochastic compositional minimization problem proposed in Section 5.4, merit further investigation. Since our proposal relies on computing projections efficiently, approaches for reducing the computational effort spent on projections, such as random constraint projection techniques, could be explored. Additionally, extensions that incorporate deterministic nonconvex constraints in a more natural fashion provide an avenue for future work. The projected stochastic subgradient method of Davis and Drusvyatskiy [19] has recently been extended to the non-Euclidean case [64], which could accelerate convergence of our proposal in practice. Finally, because stochastic approximation algorithms are an active area of research, several auxiliary techniques, such as adaptive step sizes, parallelization, acceleration, etc., may be determined to be applicable (and practically useful) to our setting.

Figure 4: Comparison of the efficient frontiers for Case study 3.



Acknowledgments

R.K. thanks Rui Chen, Eli Towle, and Clément Royer for helpful discussions.

References

- [1] Lukáš Adam and Martin Branda. Nonlinear chance constrained problems: optimality conditions, regularization and solvers. *Journal of Optimization Theory and Applications*, 170(2):419–436, 2016.
- [2] Lukáš Adam and Martin Branda. Machine learning approach to chance-constrained problems: An algorithm based on the stochastic gradient descent. http://www.optimization-online.org/DB_HTML/2018/12/6983.html (Last accessed: December 17, 2018), 2018.
- [3] Lukáš Adam, Martin Branda, Holger Heitsch, and René Henrion. Solving joint chance constrained problems using regularization and Benders decomposition. *Annals of Operations Research*, pages 1–27, 2018. doi: 10.1007/s10479-018-3091-9.
- [4] Patrick R Amestoy, Iain S Duff, Jean-Yves L’Excellent, and Jacko Koster. MUMPS: a general purpose distributed memory sparse solver. In *International Workshop on Applied Parallel Computing*, pages 121–130. Springer, 2000.
- [5] Aharon Ben-Tal, Laurent El Ghaoui, and Arkadi Nemirovski. *Robust optimization*. Princeton University Press, 2009.
- [6] Jacques F Benders. Partitioning procedures for solving mixed-variables programming problems. *Numerische Mathematik*, 4(1):238–252, 1962.
- [7] Jeff Bezanson, Alan Edelman, Stefan Karpinski, and Viral B Shah. Julia: A fresh approach to numerical computing. *SIAM review*, 59(1):65–98, 2017.
- [8] Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. *SIAM Review*, 60(2):223–311, 2018.

- [9] Giuseppe Calafiore and Marco C Campi. Uncertain convex programs: randomized solutions and confidence levels. *Mathematical Programming*, 102(1):25–46, 2005.
- [10] Giuseppe C Calafiore, Fabrizio Dabbene, and Roberto Tempo. Research on probabilistic methods for control system design. *Automatica*, 47(7):1279–1293, 2011.
- [11] Marco C Campi and Simone Garatti. A sampling-and-discarding approach to chance-constrained optimization: feasibility and optimality. *Journal of Optimization Theory and Applications*, 148(2):257–280, 2011.
- [12] Yankai Cao and Victor Zavala. A sigmoidal approximation for chance-constrained nonlinear programs. http://www.optimization-online.org/DB_FILE/2017/10/6236.pdf (Last accessed: December 17, 2018), 2017.
- [13] Yair Censor, Wei Chen, Patrick L Combettes, Ran Davidi, and Gabor T Herman. On the effectiveness of projection methods for convex feasibility problems with linear inequality constraints. *Computational Optimization and Applications*, 51(3):1065–1088, 2012.
- [14] Abraham Charnes, William W Cooper, and Gifford H Symonds. Cost horizons and certainty equivalents: an approach to stochastic programming of heating oil. *Management Science*, 4(3):235–263, 1958.
- [15] Wenqing Chen, Melvyn Sim, Jie Sun, and Chung-Piaw Teo. From CVaR to uncertainty set: Implications in joint chance-constrained optimization. *Operations Research*, 58(2):470–485, 2010.
- [16] Frank H Clarke. *Optimization and nonsmooth analysis*, volume 5. SIAM, 1990.
- [17] Laurent Condat. Fast projection onto the simplex and the ℓ_1 ball. *Mathematical Programming*, 158(1-2): 575–585, 2016.
- [18] Frank E Curtis, Andreas Wachter, and Victor M Zavala. A sequential algorithm for solving nonlinear optimization problems with chance constraints. *SIAM Journal on Optimization*, 28(1):930–958, 2018.
- [19] Damek Davis and Dmitriy Drusvyatskiy. Stochastic subgradient method converges at the rate $O(k^{-1/4})$ on weakly convex functions. *arXiv preprint arXiv:1802.02988*, 2018.
- [20] Damek Davis, Dmitriy Drusvyatskiy, Sham Kakade, and Jason D Lee. Stochastic subgradient method converges on tame functions. *arXiv preprint arXiv:1804.07795*, 2018.
- [21] Darinka Dentcheva and Gabriela Martinez. Regularization methods for optimization problems with probabilistic constraints. *Mathematical Programming*, 138(1-2):223–251, 2013.
- [22] Dmitriy Drusvyatskiy and Courtney Paquette. Efficiency of minimizing compositions of convex functions and smooth maps. *Mathematical Programming*, pages 1–56, 2018. doi: 10.1007/s10107-018-1311-3.
- [23] John C Duchi and Feng Ruan. Stochastic methods for composite and weakly convex optimization problems. *SIAM Journal on Optimization*, 28(4):3229–3259, 2018.
- [24] Iain Dunning, Joey Huchette, and Miles Lubin. JuMP: A modeling language for mathematical optimization. *SIAM Review*, 59(2):295–320, 2017.
- [25] Abebe Geletu, Armin Hoffmann, Michael Kloppel, and Pu Li. An inner-outer approximation approach to chance constrained optimization. *SIAM Journal on Optimization*, 27(3):1834–1857, 2017.
- [26] Saeed Ghadimi and Guanghui Lan. Accelerated gradient methods for nonconvex nonlinear and stochastic programming. *Mathematical Programming*, 156(1-2):59–99, 2016.
- [27] Saeed Ghadimi, Guanghui Lan, and Hongchao Zhang. Mini-batch stochastic approximation methods for nonconvex stochastic composite optimization. *Mathematical Programming*, 155(1-2):267–305, 2016.
- [28] Ambros Gleixner, Michael Bastubbe, Leon Eifler, Tristan Gally, Gerald Gamrath, Robert Lion Gottwald, Gregor Hendel, Christopher Hojny, Thorsten Koch, Marco E. Lübbecke, Stephen J. Maher, Matthias Miltenberger, Benjamin Müller, Marc E. Pfetsch, Christian Puchert, Daniel Rehfeldt, Franziska Schlösser, Christoph Schubert, Felipe Serrano, Yuji Shinano, Jan Merlin Viernickel, Matthias Walter, Fabian Wegscheider, Jonas T. Witt, and Jakob Witzig. The SCIP Optimization Suite 6.0. Technical report, Optimization Online, July 2018. URL http://www.optimization-online.org/DB_HTML/2018/07/6692.html.

- [29] Claudia Gotzes, Holger Heitsch, René Henrion, and Rüdiger Schultz. On the quantification of nomination feasibility in stationary gas networks with random load. *Mathematical Methods of Operations Research*, 84(2):427–457, 2016.
- [30] Gurobi Optimization LLC. Gurobi Optimizer Reference Manual, 2018. URL <http://www.gurobi.com>.
- [31] René Henrion. A critical note on empirical (sample average, Monte Carlo) approximation of solutions to chance constrained programs. In *IFIP Conference on System Modeling and Optimization*, pages 25–37. Springer, 2011.
- [32] L Jeff Hong, Yi Yang, and Liwei Zhang. Sequential convex approximations to joint chance constrained programs: A Monte Carlo approach. *Operations Research*, 59(3):617–630, 2011.
- [33] AM Jasour, Necdet S Aybat, and Constantino M Lagoa. Semidefinite programming for chance constrained optimization over semialgebraic sets. *SIAM Journal on Optimization*, 25(3):1411–1440, 2015.
- [34] Chi Jin, Lydia T Liu, Rong Ge, and Michael I Jordan. On the local minima of the empirical risk. *arXiv preprint arXiv:1803.09357*, 2018.
- [35] Constantino M Lagoa, Xiang Li, and Mario Sznaier. Probabilistically constrained linear programs and risk-adjusted controller design. *SIAM Journal on Optimization*, 15(3):938–951, 2005.
- [36] Guanghui Lan and Zhiqiang Zhou. Algorithms for stochastic optimization with expectation constraints. *arXiv preprint arXiv:1604.03887*, 2016.
- [37] Pu Li, Harvey Arellano-Garcia, and Günter Wozny. Chance constrained programming approach to process optimization under uncertainty. *Computers & Chemical Engineering*, 32(1-2):25–45, 2008.
- [38] James Luedtke. A branch-and-cut decomposition algorithm for solving chance-constrained mathematical programs with finite support. *Mathematical Programming*, 146(1-2):219–244, 2014.
- [39] James Luedtke and Shabbir Ahmed. A sample approximation approach for optimization with probabilistic constraints. *SIAM Journal on Optimization*, 19(2):674–699, 2008.
- [40] Bruce L Miller and Harvey M Wagner. Chance constrained programming with joint constraints. *Operations Research*, 13(6):930–945, 1965.
- [41] Arkadi Nemirovski and Alexander Shapiro. Convex approximations of chance constrained programs. *SIAM Journal on Optimization*, 17(4):969–996, 2006.
- [42] Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609, 2009.
- [43] BK Pagnoncelli, Shapiro Ahmed, and Alexander Shapiro. Sample average approximation method for chance constrained programming: theory and applications. *Journal of Optimization Theory and Applications*, 142(2):399–416, 2009.
- [44] András Prékopa. On probabilistic constrained programming. In *Proceedings of the Princeton symposium on mathematical programming*, pages 113–138. Princeton, NJ, 1970.
- [45] András Prékopa. *Stochastic programming*, volume 324. Springer Science & Business Media, 1995.
- [46] András Prékopa and Tamás Szántai. *On optimal regulation of a storage level with application to the water level regulation of a lake*, volume I of *Studies in Applied Stochastic Programming*. 1979.
- [47] Tara Rengarajan and David P Morton. Estimating the efficient frontier of a probabilistic bicriteria model. In *Winter Simulation Conference*, pages 494–504, 2009.
- [48] R Tyrrell Rockafellar and Stanislav Uryasev. Optimization of conditional value-at-risk. *Journal of Risk*, 2:21–42, 2000.
- [49] R Tyrrell Rockafellar and Roger J-B Wets. *Variational analysis*, volume 317. Springer Science & Business Media, 2009.
- [50] JO Royset and E Polak. Reliability-based optimal design using sample average approximations. *Probabilistic Engineering Mechanics*, 19(4):331–343, 2004.

- [51] Harold Ruben. Probability content of regions under spherical normal distributions, IV: The distribution of homogeneous and non-homogeneous quadratic functions of normal variables. *The Annals of Mathematical Statistics*, 33(2):542–570, 1962.
- [52] F Shan, XT Xiao, and LW Zhang. Convergence analysis on a smoothing approach to joint chance constrained programs. *Optimization*, 65(12):2171–2193, 2016.
- [53] Feng Shan, Liwei Zhang, and Xiantao Xiao. A smoothing function approach to joint chance-constrained programs. *Journal of Optimization Theory and Applications*, 163(1):181–199, 2014.
- [54] Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczyński. *Lectures on stochastic programming: modeling and theory*. SIAM, 2009.
- [55] Wim van Ackooij and René Henrion. Gradient formulae for nonlinear probabilistic constraints with Gaussian and Gaussian-like distributions. *SIAM Journal on Optimization*, 24(4):1864–1889, 2014.
- [56] Wim van Ackooij and René Henrion. (Sub-)Gradient formulae for probability functions of random inequality systems under Gaussian distribution. *SIAM/ASA Journal on Uncertainty Quantification*, 5(1):63–87, 2017.
- [57] Wim van Ackooij, Antonio Frangioni, and Welington de Oliveira. Inexact stabilized Benders decomposition approaches with application to chance-constrained problems with finite support. *Computational Optimization and Applications*, 65(3):637–669, 2016.
- [58] Wim van Ackooij, V Berge, Welington de Oliveira, and C Sagastizábal. Probabilistic optimization via approximate p-efficient points and bundle methods. *Computers & Operations Research*, 77:177–193, 2017.
- [59] Andreas Wächter and Lorenz T Biegler. On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming. *Mathematical programming*, 106(1):25–57, 2006.
- [60] Mengdi Wang and Dimitri P Bertsekas. Stochastic first-order methods with random constraint projection. *SIAM Journal on Optimization*, 26(1):681–717, 2016.
- [61] Mengdi Wang, Ethan X Fang, and Han Liu. Stochastic compositional gradient descent: algorithms for minimizing compositions of expected-value functions. *Mathematical Programming*, 161(1-2):419–449, 2017.
- [62] Xiao Wang, Shiqian Ma, and Ya-xiang Yuan. Penalty methods with stochastic approximation for stochastic nonlinear programming. *Mathematics of Computation*, 86(306):1793–1820, 2017.
- [63] Hui Zhang and Pu Li. Chance constrained programming for optimal power flow under uncertainty. *IEEE Transactions on Power Systems*, 26(4):2417–2424, 2011.
- [64] Siqi Zhang and Niao He. On the convergence rate of stochastic mirror descent for nonsmooth nonconvex optimization. *arXiv preprint arXiv:1806.04781*, 2018.

A Algorithm outlines

A.1 Implementation of the scenario approximation algorithm

Algorithm 6 details our implementation of the tuned scenario approximation algorithm. We use $M = 50$ iterations and $R = 20$ replicates for our computational experiments. Settings for sample sizes N_i , constraints added per iteration N_c , and number of samples N_{MC} are problem dependent, and provided below. Line 19 in Algorithm 6, which estimates the risk level of a candidate solution, may be replaced, if possible, by one that computes the true analytical risk level or a numerical estimate of it [55].

- $N_i = \lceil 10^{a_i} \rceil$, $i = 1, \dots, 50$, where: for Case study 1: $a_i = 1 + \frac{5}{49}(i - 1)$; for Case study 2: $a_i = 1 + \frac{\log_{10}(50000) - 1}{49}(i - 1)$; for Case study 3: $a_i = 1 + \frac{4}{49}(i - 1)$; and for Case study 4: $a_i = 1 + \frac{4}{49}(i - 1)$ (the upper bounds on N_i were determined based on memory requirements)
- N_c : Case study 1: 1000; Case study 2: 10; Case study 3: 5; and Case study 4: 10 (these values were tuned for good performance)
- N_{MC} : See Section 6.1 of the paper (we use the same Monte Carlo samples that we used for our proposed method to estimate risk levels).

Algorithm 6 Scenario approximation for approximating the efficient frontier of (CCP)

1: **Input:** Number of major iterations M , distinct sample sizes N_i , $i = 1, \dots, M$, number of replicates per major iteration R , maximum number of scenario constraints added per iteration N_c , number of samples to estimate risk levels N_{MC} , and initial guess $\hat{x} \in X$.

2: **Output:** Pairs $(\bar{\nu}^{i,r}, \bar{\alpha}^{i,r})$, $i = 1, \dots, M$, $r = 1, \dots, R$, of objective values and risk levels that can be used to approximate the efficient frontier, and corresponding sequence of solutions $\{\bar{x}^{i,r}\}$.

3: **Preparation:** Draw a (fixed) sample $\{\bar{\xi}^l\}_{l=1}^{N_{MC}}$ from \mathcal{P} for estimating risk levels of candidate solutions.

4: **for** major iteration $i = 1$ to M **do**

5: **for** replicate $r = 1$ to R **do**

6: Sample: Draw i.i.d. sample $\{\xi^{l,r}\}_{l=1}^{N_i}$ from \mathcal{P} .

7: Initialize: Starting point $x^0 := \hat{x}$, list of scenarios enforced for each chance constraint $\mathcal{I}_k^0 = \emptyset$, $k = 1, \dots, m$, and counter $q = 0$.

8: **repeat**

9: Update $q \leftarrow q + 1$, number of scenario constraints violated $V \leftarrow 0$, $\mathcal{I}_k^q = \mathcal{I}_k^{q-1}$, $\forall k$.

10: Solve the following scenario approximation problem (locally) using the initial guess x^{q-1} to obtain a solution x^q and corresponding optimal objective ν^q :

$$\begin{aligned} \min_{x \in X} \quad & c^T x \\ \text{s.t.} \quad & g_k(x, \xi^{l,r}) \leq 0, \quad \forall l \in \mathcal{I}_k^{q-1}. \end{aligned}$$

11: **for** $k = 1$ to m **do**

12: Evaluate k^{th} random constraint g_k at x^q for each of the scenarios $l = 1, \dots, N_i$.

13: Sort the values $\{g_k(x^q, \xi^{l,r})\}_{l=1}^{N_i}$ in decreasing order.

14: Increment V by the number of scenario constraints satisfying $g_k(x^q, \xi^{l,r}) > 0$.

15: Add at most N_c of the most violated scenario indices l to \mathcal{I}_k^q .

16: **end for**

17: **until** $V = 0$

18: Set $\hat{x} = x^{i,r} = x^q$ and $\bar{\nu}^{i,r} = \nu^q$ to their converged values.

19: Estimate risk level $\bar{\alpha}^{i,r}$ of candidate solution \hat{x} using the sample $\{\bar{\xi}^l\}_{l=1}^{N_{MC}}$.

20: **end for**

21: **end for**

A.2 Solving (CCP) for a fixed risk level

Algorithm 7 adapts Algorithm 1 to solve (CCP) for a given risk level $\hat{\alpha} \in (0, 1)$. An initial point \bar{x}^0 and an upper bound on the optimal objective value ν_{up} can be obtained in a manner similar to Algorithm 1, whereas a lower bound on the optimal objective value ν_{low} can be obtained either using lower bounding techniques (see [39, 41]), or by trial and error. Note that the sequence of approximations in line 9 of Algorithm 7 need not be solved until termination if Algorithm 2 determines that $\bar{\alpha}^i < \hat{\alpha}$ before its termination criteria have been satisfied.

Algorithm 7 Solving (CCP) for a fixed risk level

- 1: **Input:** target risk level $\hat{\alpha} \in (0, 1)$, ‘guaranteed’ lower and upper bounds on the optimal objective value ν_{low} and ν_{up} with $\nu_{low} \leq \nu_{up}$, and initial point $\bar{x}^0 \in X$.
 - 2: **Set algorithmic parameters:** in addition to line 2 of Algorithm 1, let $\nu_{tol} > 0$ denote an optimality tolerance.
 - 3: **Output:** approximate optimal objective value $\hat{\nu}$ of (CCP) for a risk level of $\hat{\alpha}$.
 - 4: **Preprocessing:** let $\bar{\nu}^0 = \frac{1}{2}(\nu_{low} + \nu_{up})$, and determine smoothing parameters $\{\tau_{k,j}\}_{k=1}^K$ scaled at $\text{proj}(\bar{x}^0, X_{\bar{\nu}^0})$ using Algorithm 3 and an initial sequence of step lengths $\{\bar{\gamma}_k\}_{k=1}^K$ for the corresponding sequence of approximating problems (APP) using Algorithm 4.
 - 5: **Optimization Phase:**
 - 6: Initialize index $i = 0$.
 - 7: **repeat**
 - 8: Update iteration count $i \leftarrow i + 1$ and set objective bound $\bar{\nu}^i = \frac{1}{2}(\nu_{low} + \nu_{up})$.
 - 9: Obtain $(\bar{\alpha}^i, \bar{x}^i)$ by solving sequence of approximations (APP) using Algorithm 2 with the above algorithmic parameter settings and $\text{proj}(\bar{x}^{i-1}, X_{\bar{\nu}^i})$ as the initial guess.
 - 10: **if** $\bar{\alpha}^i \geq \hat{\alpha}$ **then** set $\nu_{low} \leftarrow \bar{\nu}^i$; **else** set $\nu_{up} \leftarrow \bar{\nu}^i$.
 - 11: **until** $\nu_{low} \geq \nu_{up} - \nu_{tol}$
 - 12: Set $\hat{\nu} = \nu_{up}$.
-

A.3 Implementation of the sigmoidal approximation algorithm

We used the following ‘tuned’ settings to solve each iteration of the sigmoidal approximation problem (see Section 4 of Cao and Zavala [12]) using IPOPT: `tol` = 10^{-4} , `max_iter` = 10000, `hessian_approximation` = `limited_memory`, `jac_c_constant`=yes, and `max_cpu_time`=3600 seconds. We terminated the loop of Algorithm SigVar-Alg of Cao and Zavala [12] when the objective improved by less than 0.01% relative to the previous iteration. In what follows, we use the notation of Algorithm SigVar-Alg of Cao and Zavala [12]. For our ‘adaptive γ ’ setting, we use the proposal of Cao and Zavala [12] to specify γ when the solution of the CVaR problem corresponds to $t_c(\alpha) < 0$. When $t_c(\alpha) = 0$ at the CVaR solution returned by Gurobi, we try to estimate a good value of $t_c(\alpha)$ by looking at the true distribution of $g(x_c(\alpha), \xi)$.

B Proofs

B.1 Proof of Proposition 1

Define $h_k : \mathbb{R}^n \times \mathbb{R}^d \rightarrow \mathbb{R}$ by $h_k(x, \xi) := \max[\phi_k(g(x, \xi))]$ for each $k \in \mathbb{N}$, and note that h_k is continuous by virtue of Assumption (4a). Additionally, Assumption (4c) implies that $|h_k(x, \xi)| \leq M_\phi$ for each $x \in \mathbb{R}^n, \xi \in \mathbb{R}^d$, and $k \in \mathbb{N}$. By noting that for each $x \in X_\nu$,

$$\lim_{k \rightarrow \infty} h_k(x, \xi) = \max[\mathbb{1}[g(x, \xi)]]$$

for a.e. $\xi \in \Xi$ due to Assumptions 2 and 4 (or, simply, the strong form of Assumption 4), we obtain the stated result by Lebesgue’s dominated convergence theorem [54, Theorem 7.31]. \square

B.2 Proof of Proposition 2

Define $h_k : \mathbb{R}^n \rightarrow \mathbb{R}$ and $h : \mathbb{R}^n \rightarrow \mathbb{R}$ by $h_k(x) := \hat{p}_k(x) + I_{X_\nu}(x)$ and $h(x) := p(x) + I_{X_\nu}(x)$, respectively. From Proposition 7.2 of Rockafellar and Wets [49], we have that $\{h_k\}$ epi-converges to h if and only if

at each $x \in \mathbb{R}^n$

$$\begin{aligned} \liminf_{k \rightarrow \infty} h_k(x_k) &\geq h(x), \text{ for every sequence } \{x_k\} \rightarrow x, \text{ and} \\ \limsup_{k \rightarrow \infty} h_k(x_k) &\leq h(x), \text{ for some sequence } \{x_k\} \rightarrow x. \end{aligned}$$

Consider the constant sequence with $x_k = x$ for each $k \in \mathbb{N}$. We have

$$\limsup_{k \rightarrow \infty} h_k(x_k) = I_{X_\nu}(x) + \limsup_{k \rightarrow \infty} \hat{p}_k(x) = I_{X_\nu}(x) + p(x) = h(x)$$

as a result of Proposition 1, which establishes the latter inequality.

To see the former inequality, consider an arbitrary sequence $\{x_k\}$ in \mathbb{R}^n converging to $x \in \mathbb{R}^n$. By noting that the characteristic function $I_{X_\nu}(\cdot)$ is lower semicontinuous (since X_ν is closed) and

$$\liminf_{k \rightarrow \infty} h_k(x_k) \geq \liminf_{k \rightarrow \infty} \hat{p}_k(x_k) + \liminf_{k \rightarrow \infty} I_{X_\nu}(x_k),$$

it suffices to show $\liminf_k \hat{p}_k(x_k) \geq p(x)$ to establish $\{h_k\}$ epi-converges to h . In fact, it suffices to show the above inequality holds for any $x \in X_\nu$ since the former inequality holds trivially for $x \notin X_\nu$.

Define $q_k : \mathbb{R}^d \rightarrow \mathbb{R}$ by $q_k(\xi) := \max [\phi_k(g(x_k, \xi))]$, and note that q_k is \mathbb{P} -integrable by virtue of Assumption 4. By Fatou's lemma [54, Theorem 7.30], we have that

$$\liminf_k \mathbb{E}[q_k(\xi)] = \liminf_k \hat{p}_k(x_k) \geq \mathbb{E}[q(\xi)],$$

where $q : \mathbb{R}^d \rightarrow \mathbb{R}$ is defined as

$$q(\xi) := \liminf_k q_k(\xi) = \liminf_k \max [\phi_k(g(x_k, \xi))].$$

Therefore, it suffices to show that $q(\xi) \geq \max [\mathbb{1}[g(x, \xi)]]$ a.e. for each $x \in X_\nu$ to establish the former inequality (we note that this result holds sans Assumption 2 if the strong form of Assumption 4 is made).

Fix $\xi \in \Xi$. Let $l \in \{1, \dots, m\}$ denote an index at which the maximum in $\max [\mathbb{1}[g(x, \xi)]]$ is attained, i.e., $\max [\mathbb{1}[g(x, \xi)]] = \mathbb{1}[g_l(x, \xi)]$. Consider first the case when $g_l(x, \xi) > 0$. By the continuity of the functions g , for any $\varepsilon > 0$ there exists $N_\varepsilon \in \mathbb{N}$ such that $j \geq N_\varepsilon \implies g(x_j, \xi) > g(x, \xi) - \varepsilon$. Therefore

$$q(\xi) = \liminf_k \max [\phi_k(g(x_k, \xi))] \geq \liminf_k \max [\phi_k(g(x, \xi) - \varepsilon)] \geq \liminf_k \phi_{k,l}(g_l(x, \xi) - \varepsilon),$$

where the first inequality follows from Assumption (4b). Choosing $\varepsilon < g_l(x, \xi)$ yields

$$q(\xi) = \liminf_k \max [\phi_k(g(x_k, \xi))] \geq \liminf_k \phi_{k,l}(g_l(x, \xi) - \varepsilon) = 1 = \max [\mathbb{1}[g(x, \xi)]]$$

by virtue of Assumption (4d). The case when $g_l(x, \xi) < 0$ follows more directly since

$$q(\xi) = \liminf_k \max [\phi_k(g(x_k, \xi))] \geq \liminf_k \phi_{k,l}(g_l(x, \xi) - \varepsilon) = \mathbb{1}[g_l(x, \xi) - \varepsilon] = 0 = \max [\mathbb{1}[g(x, \xi)]],$$

where the first equality follows from Assumption 4. Since $\max [g(x, \xi)] \neq 0$ a.e. $\forall x \in X_\nu$ by Assumption 2, we have that $q(\xi) \geq \max [\mathbb{1}[g(x, \xi)]]$ a.e. for each $x \in X_\nu$, which concludes the proof. \square

B.3 Proof of Proposition 4

Since x^* is a strict local minimizer of (1), there exists $\delta > 0$ such that $p(x) > p(x^*)$, $\forall x \in X_\nu \cap \text{cl}(B_\delta(x^*))$. Since $X_\nu \cap \text{cl}(B_\delta(x^*))$ is compact, $\min_{x \in X_\nu \cap \text{cl}(B_\delta(x^*))} \hat{p}_k(x)$ has a global minimum, say \hat{x}_k , due to Assumption (4a). Furthermore, $\{\hat{x}_k\}$ has a convergent subsequence in $X_\nu \cap \text{cl}(B_\delta(x^*))$.

Assume without loss of generality that $\{\hat{x}_k\}$ itself converges to $\hat{x} \in X_\nu \cap \text{cl}(B_\delta(x^*))$. Applying Theorem 1 to the above sequence of restricted minimization problems yields the conclusion that \hat{x} is a global minimizer of p on $X_\nu \cap \text{cl}(B_\delta(x^*))$, i.e., $p(x) \geq p(\hat{x})$, $\forall x \in X_\nu \cap \text{cl}(B_\delta(x^*))$. Since x^* is a strict local minimizer of (1), this implies $\hat{x} = x^*$. Since $\{\hat{x}_k\} \rightarrow x^*$, this implies that x_k belongs to the open ball $B_\delta(x^*)$ for k sufficiently large. Therefore, \hat{x}_k is a local minimizer of $\min_{x \in X_\nu} \hat{p}_k(x)$ for k sufficiently large since it is a global minimizer of the above problem on $B_\delta(x^*)$. \square

B.4 Proof of Lemma 2

1. For any $x, y \in X_\nu$, $j \in \{1, \dots, m\}$, and $\xi \in \Xi$, we have

$$|\phi_{k,j}(g_j(y, \xi)) - \phi_{k,j}(g_j(x, \xi))| \leq M'_{\phi,k,j} |g_j(y, \xi) - g_j(x, \xi)| \leq M'_{\phi,k,j} L_{g,j}(\xi) \|y - x\|,$$

where the first step follows from Assumptions (4a) and (5a) and the mean-value theorem, and the second step follows from Assumption (3a).

2. For any $x, y \in X_\nu$, we have

$$\begin{aligned} |\hat{p}_k(y) - \hat{p}_k(x)| &= |\mathbb{E}[\max[\phi_k(g(y, \xi))] - \max[\phi_k(g(x, \xi))]]| \\ &\leq \mathbb{E}[\max[|\phi_k(g(y, \xi)) - \phi_k(g(x, \xi))|]] \\ &\leq \mathbb{E}\left[\max_j \left[M'_{\phi,k,j} L_{g,j}(\xi)\right]\right] \|y - x\|. \quad \square \end{aligned}$$

B.5 Proof of Proposition 5

From the corollary to Proposition 2.2.1 of Clarke [16], we have that $\phi_k(g(\cdot, \xi))$ is strictly differentiable. Theorem 2.3.9 of Clarke [16] then implies that $\max[\phi_k(g(\cdot, \xi))]$ is Clarke regular. The stronger assumptions of Theorem 2.7.2 of Clarke [16] are satisfied because of Assumption 4, Lemma 2, and the regularity of $\max[\phi_k(g(\cdot, \xi))]$, which then yields

$$\partial \hat{p}_k(x) = \partial \mathbb{E}[\max[\phi_k(g(x, \xi))] = \mathbb{E}[\partial_x \max[\phi_k(g(x, \xi))]].$$

Noting that $\phi_k(g(\cdot, \xi))$ is Clarke regular from Proposition 2.3.6 of Clarke [16], the stated equality then follows from Proposition 2.3.12 of Clarke [16]. \square

B.6 Proof of Proposition 6

We first sketch an outline of the proof (note that our proof will follow a different order for reasons that will become evident). We will establish that (see Chapters 4 and 5 of Rockafellar and Wets [49] for definitions of technical terms)

$$\limsup_{\substack{x \rightarrow \bar{x} \\ k \rightarrow \infty}} \partial \hat{p}_k(x) + N_{X_\nu}(x) = \limsup_{\substack{x \rightarrow \bar{x} \\ k \rightarrow \infty}} \partial \hat{p}_k(x) + \limsup_{\substack{x \rightarrow \bar{x} \\ k \rightarrow \infty}} N_{X_\nu}(x). \quad (3)$$

Then, because of the outer semicontinuity of the normal cone mapping, it suffices to prove that the outer limit of $\partial \hat{p}_k(x)$ is a subset of $\{\nabla p(\bar{x})\}$. To demonstrate this, we will first show that $\partial \hat{p}_k(x) = -\partial \int_{-\infty}^{\infty} F(x, \eta) d\bar{\phi}_k(\eta) d\eta$, where F is the cumulative distribution function of $\max[\bar{g}(x, \xi)]$. Then, we will split this integral into two parts - one accounting for tail contributions (which we will show vanishes when we take the outer limit), and the other accounting for the contributions of F near $(x, \eta) = (\bar{x}, 0)$ that we will show satisfies the desired outer semicontinuity property.

Since $\hat{p}_k(x)$ can be rewritten as $\hat{p}_k(x) = \mathbb{E}[\bar{\phi}_k(\max[\bar{g}(x, \xi)])]$ by Assumption 6, we have

$$\begin{aligned} \hat{p}_k(x) &= \mathbb{E}[\bar{\phi}_k(\max[\bar{g}(x, \xi)])] = \int_{-\infty}^{+\infty} \bar{\phi}_k(\max[\bar{g}(x, \xi)]) d\mathbb{P} \\ &= \int_{-\infty}^{+\infty} \bar{\phi}_k(\eta) dF(x, \eta) \\ &= \lim_{\eta \rightarrow +\infty} \bar{\phi}_k(\eta) - \int_{-\infty}^{+\infty} F(x, \eta) d\bar{\phi}_k(\eta) d\eta, \end{aligned}$$

where the second line is to be interpreted as a Lebesgue-Stieljes integral, and the final step follows by integrating by parts and Assumption 4. This yields (see Proposition 5 and Assumption 6 for the existence of these quantities)

$$\partial \hat{p}_k(x) \subset -\partial \int_{|\eta| \geq \varepsilon_k} F(x, \eta) d\bar{\phi}_k(\eta) d\eta - \partial \int_{-\varepsilon_k}^{\varepsilon_k} F(x, \eta) d\bar{\phi}_k(\eta) d\eta \quad (4)$$

by the properties of the Clarke generalized gradient, where $\{\varepsilon_k\}$ is defined in Assumption 6.

Let $\{x_k\}$ be any sequence in X_ν converging to $\bar{x} \in X_\nu$. Suppose $v_k \in -\partial \int_{|\eta| \geq \varepsilon_k} F(x_k, \eta) d\bar{\phi}_k(\eta) d\eta$ with $v_k \rightarrow v \in \mathbb{R}^n$. We would like to show that $v = 0$. Note that by an abuse of notation (where we actually take norms of the integrable selections that define v_k)

$$\begin{aligned} \lim_{k \rightarrow \infty} \|v_k\| &= \lim_{k \rightarrow \infty} \left\| -\partial \int_{|\eta| \geq \varepsilon_k} F(x_k, \eta) d\bar{\phi}_k(\eta) d\eta \right\| = \lim_{k \rightarrow \infty} \left\| \int_{|\eta| \geq \varepsilon_k} \partial_x F(x_k, \eta) d\bar{\phi}_k(\eta) d\eta \right\| \\ &\leq \lim_{k \rightarrow \infty} \int_{|\eta| \geq \varepsilon_k} \|\partial_x F(x_k, \eta)\| d\bar{\phi}_k(\eta) d\eta \\ &\leq \lim_{k \rightarrow \infty} \int_{|\eta| \geq \varepsilon_k} L_F(\eta) d\bar{\phi}_k(\eta) d\eta = 0, \end{aligned} \quad (5)$$

where the first step follows from Theorem 2.7.2 of Clarke [16] (whose assumptions are satisfied by virtue of Assumption (6b)), the second inequality follows from Proposition 2.1.2 of Clarke [16], and the final equality follows from Assumption (6c).

The above arguments establish that $\limsup_{\substack{x \rightarrow \bar{x} \\ k \rightarrow \infty}} -\partial \int_{|\eta| \geq \varepsilon_k} F(x, \eta) d\bar{\phi}_k(\eta) d\eta = \{0\}$. Consequently, we have $\limsup_{\substack{x \rightarrow \bar{x} \\ k \rightarrow \infty}} \partial \hat{p}_k(x) \subset \limsup_{\substack{x \rightarrow \bar{x} \\ k \rightarrow \infty}} -\partial \int_{-\varepsilon_k}^{\varepsilon_k} F(x, \eta) d\bar{\phi}_k(\eta) d\eta$ from Equation (4). We now consider the outer limit $\limsup_{\substack{x \rightarrow \bar{x} \\ k \rightarrow \infty}} -\partial \int_{-\varepsilon_k}^{\varepsilon_k} F(x, \eta) d\bar{\phi}_k(\eta) d\eta$. Suppose $w_k \in -\partial \int_{-\varepsilon_k}^{\varepsilon_k} F(x_k, \eta) d\bar{\phi}_k(\eta) d\eta$ with $w_k \rightarrow w \in \mathbb{R}^n$. We wish to show that $w \in \{\nabla p(\bar{x})\}$. Invoking Theorem 2.7.2 of Clarke [16] once again, we have that

$$\begin{aligned} \limsup_{k \rightarrow \infty} -\partial \int_{-\varepsilon_k}^{\varepsilon_k} F(x_k, \eta) d\bar{\phi}_k(\eta) d\eta &\subset \limsup_{k \rightarrow \infty} -\int_{-\varepsilon_k}^{\varepsilon_k} \nabla_x F(x_k, \eta) d\bar{\phi}_k(\eta) d\eta \\ &= \limsup_{k \rightarrow \infty} -\left(\int_{-\varepsilon_k}^{\varepsilon_k} \nabla_x F(x_k, \eta) d\hat{\phi}_k(\eta) d\eta \right) \left(\int_{-\varepsilon_k}^{\varepsilon_k} d\bar{\phi}_k(z) dz \right) \\ &= \left(\limsup_{k \rightarrow \infty} -\int_{-\varepsilon_k}^{\varepsilon_k} \nabla_x F(x_k, \eta) d\hat{\phi}_k(\eta) d\eta \right) \left(\lim_{k \rightarrow \infty} \int_{-\varepsilon_k}^{\varepsilon_k} d\bar{\phi}_k(z) dz \right) \\ &= \limsup_{k \rightarrow \infty} -\int_{-\varepsilon_k}^{\varepsilon_k} \nabla_x F(x_k, \eta) d\hat{\phi}_k(\eta) d\eta, \end{aligned} \quad (6)$$

where for each $k \in \mathbb{N}$ large enough, $\hat{\phi}_k : \mathbb{R} \rightarrow \mathbb{R}$ is defined as $\hat{\phi}_k(y) = \frac{\bar{\phi}_k(y)}{\int_{-\varepsilon_k}^{\varepsilon_k} d\bar{\phi}_k(z) dz}$, the first step follows (by an abuse of notation) from the fact that $\varepsilon_k < \theta$ for k large enough (see Assumption (6b)), and the third and fourth steps follow from Assumption (6c). Noting that

$$\int_{-\varepsilon_k}^{\varepsilon_k} \nabla_x F(x_k, \eta) d\hat{\phi}_k(\eta) d\eta = \left[\frac{\partial F}{\partial x_1}(x_k, \omega_{1,k}) \cdots \frac{\partial F}{\partial x_n}(x_k, \omega_{n,k}) \right]^T \int_{-\varepsilon_k}^{\varepsilon_k} d\hat{\phi}_k(\eta) d\eta = \begin{bmatrix} \frac{\partial F}{\partial x_1}(x_k, \omega_{1,k}) \\ \vdots \\ \frac{\partial F}{\partial x_n}(x_k, \omega_{n,k}) \end{bmatrix}$$

for some constants $\omega_{i,k} \in (-\varepsilon_k, \varepsilon_k)$, $i = 1, \dots, n$, by virtue of Assumption (6b) and the first mean value theorem for definite integrals, we have

$$\begin{aligned} \limsup_{k \rightarrow \infty} -\partial \int_{-\varepsilon_k}^{\varepsilon_k} F(x_k, \eta) d\bar{\phi}_k(\eta) d\eta &\subset \limsup_{k \rightarrow \infty} -\left[\frac{\partial F}{\partial x_1}(x_k, \omega_{1,k}) \cdots \frac{\partial F}{\partial x_n}(x_k, \omega_{n,k}) \right]^T \\ &= \{-\nabla_x F(\bar{x}, 0)\} = \{\nabla p(\bar{x})\}, \end{aligned}$$

where the first equality above follows from Assumption (6b), and the second equality follows from the fact that $p(x) = 1 - F(x, 0)$ for each $x \in X_\nu$. Reconciling our progress with Equation (4), we obtain

$$\limsup_{\substack{x \rightarrow \bar{x} \\ k \rightarrow \infty}} \partial \hat{p}_k(x) \subset \{\nabla p(\bar{x})\}.$$

To establish the desirable equality in Equation (3), it suffices to show that $\partial \hat{p}_k(x_k) \subset C$ for k large enough, where $C \subset \mathbb{R}^n$ is independent of k . From Equation (5), we have that the first term in the right-hand side of Equation (4) is contained in a bounded set that is independent of $k \in \mathbb{N}$. From Equation (6), we have that any element of the second term in the right-hand side of Equation (4) is bounded above in norm by $\max_{\substack{x \in \text{cl}(B_\delta(\bar{x})) \\ \eta \in [-0.5\theta, 0.5\theta]}} \|\nabla_x F(x, \eta)\|$ for k large enough for any $\delta > 0$. The above arguments in conjunction

with Equation (4) establish Equation (3). The desired result then follows from the outer semicontinuity of the normal cone mapping, see Proposition 6.6 of Rockafellar and Wets [49]. \square

B.7 Proof of Proposition 7

The bounds follow by noting that for each $y \in \mathbb{R}$:

$$d\phi_k(y) = \frac{\tau_{k,j}^{-1}}{\left[\exp(0.5\tau_{k,j}^{-1}y) + \exp(-0.5\tau_{k,j}^{-1}y)\right]^2} \leq 0.25\tau_{k,j}^{-1}, \text{ and}$$

$$|d^2\phi_{k,j}(y)| = \tau_{k,j}^{-2} \left| \frac{\exp\left(\frac{y}{2\tau_{k,j}}\right) - \exp\left(-\frac{y}{2\tau_{k,j}}\right)}{\left[\exp\left(\frac{y}{2\tau_{k,j}}\right) + \exp\left(-\frac{y}{2\tau_{k,j}}\right)\right]^3} \right| \leq 0.1\tau_{k,j}^{-2},$$

where the final inequality is obtained by maximizing the function $x \mapsto \left| \frac{x - \frac{1}{x}}{\left(x + \frac{1}{x}\right)^3} \right|$ on \mathbb{R} . \square

B.8 Proof of Proposition 9

1. Follows by using the chain rule $\nabla \phi_{k,j}(g_j(\cdot, \xi)) = d\phi_{k,j}(g_j(\cdot, \xi)) \nabla g_j(\cdot, \xi)$ and bounding the Lipschitz constant of the above product on X_ν using Assumptions (3a), (3b), and 5 and standard arguments.
2. From part one of Lemma 2 and Theorem 7.44 of Shapiro et al. [54], we have $\nabla \mathbb{E}[\phi_k(g(x, \xi))] = \mathbb{E}[\nabla_x \phi_k(g(x, \xi))]$, $\forall x \in X_\nu$. The stated result then follows from the first part of this proposition and the fact that the Frobenius norm of a matrix provides an upper bound on its spectral norm, where the existence of the quantity $\mathbb{E} \left[\left(\sum_{j=1}^m L_{k,j}^2(\xi) \right)^{\frac{1}{2}} \right]$ follows from the fact that $\sqrt{\sum_j L_{k,j}^2(\xi)} \leq \sum_j L_{k,j}(\xi)$, Assumptions (3a) and (3b), and Lebesgue's dominated convergence theorem. \square

C Recourse formulation

As noted in Section 1 of the paper, Problem (CCP) can also be used to model chance-constrained programs with static recourse decisions, e.g., by defining g through the solution of the following auxiliary optimization problem for each $(x, \xi) \in X \times \Xi$:

$$\begin{aligned} g(x, \xi) &:= \min_{y, \eta} \eta \\ \text{s.t. } & T(x, \xi) + W(\xi)y + \eta e \geq 0, \\ & y \in \mathbb{R}_+^{n_y}, \quad \eta \geq -1, \end{aligned}$$

where $T : \mathbb{R}^n \times \mathbb{R}^d \rightarrow \mathbb{R}^{m_r}$ is continuously differentiable, and $W : \mathbb{R}^d \rightarrow \mathbb{R}^{m_r \times n_y}$. Note that g can be recast in the form of a joint chance constraint by appealing to linear programming duality, viz.,

$$g(x, \xi) = \max_{j=1, \dots, |V(\xi)|} -v_j^T T(x, \xi) - w_j,$$

where the polytope $P_g(\xi) = \{(v, w) \in \mathbb{R}_+^{m_r} \times \mathbb{R}_+ : v^T W(\xi) \leq 0, v^T e + w = 1\}$, and $V(\xi)$ denotes the set of extreme points of $P_g(\xi)$. Reformulating the recourse constraints into the above explicit form may not be practical for our proposal since Algorithms 3 and 4 rely on stepping through each of the constraint functions g_j , $j = 1, \dots, |V(\xi)|$, and the cardinality of the set of extreme points $V(\xi)$ may

be huge. Therefore, we consider the case when the approximations \hat{p}_k are constructed using a single smoothing function, i.e., $\hat{p}_k(x) = \mathbb{E}[\phi_k(g(x, \xi))]$, where $\{\phi_k\}$ is a sequence of smooth scalar functions that approximate the step function.

Throughout this section, we assume that g that can be reformulated as $g(x, \xi) := \max_{j \in J} h_j(x, \xi)$, where J is a finite index set, and $h_j : \mathbb{R}^n \times \mathbb{R}^d \rightarrow \mathbb{R}$ are continuously differentiable functions that are Lipschitz continuous and have Lipschitz continuous gradients in the sense of Assumption 3. We let $L'_h(\xi)$ denote the Lipschitz constant of $\nabla h(\cdot, \xi)$ on X_ν with $\mathbb{E} \left[\left(L'_h(\xi) \right)^2 \right] < +\infty$. While the theory in this section is also applicable to the ordinary joint chance constrained setting, it is usually advantageous to consider individual smoothing functions for each chance constraint function in that case whenever feasible. As an aside, we mention that imposing the lower bound constraint $\eta \geq -1$ in the definition of g is not necessary (and it may in fact be advantageous to omit it) when the set

$$\left\{ z \in \mathbb{R} : z = \min_{i=1, \dots, m_r} T_i(x, \xi) + (w_i(\xi))^T y \text{ for some } (x, y, \xi) \in X \times \mathbb{R}_+^{n_y} \times \Xi \right\}$$

is bounded from above, where $w_i(\xi)$ denotes the vector corresponding to the i^{th} row of $W(\xi)$.

We first characterize the Clarke generalized gradient of the approximation \hat{p}_k under the above setup.

Proposition 11. Suppose Assumptions 4 and 5 and the above assumptions on $g(x, \xi) = \max[h(x, \xi)]$ and ϕ_k hold. Then $\partial \hat{p}_k(x) = \mathbb{E}[d\phi_k(\max[h(x, \xi)]) \times \partial_x \max[h(x, \xi)]]$.

Proof. Note that for any $x \in X$ and $\xi \in \Xi$, we have $\partial_x \phi_k(g(x, \xi)) = d\phi_k(g(x, \xi)) \partial_x \max[h(x, \xi)]$, see Theorem 2.3.9 of Clarke [16]. The claim then follows from Proposition 5 of the paper and the fact that $\hat{p}_k(x) = \mathbb{E}[\phi_k(g(x, \xi))] = \mathbb{E} \left[\max_j [\phi_k(g_j(x, \xi))] \right]$ since ϕ_k is monotonically nondecreasing on \mathbb{R} . \square

The following result establishes that the approximation \hat{p}_k continues to enjoy the weak convexity property for the above setup under mild assumptions.

Proposition 12. Suppose Assumptions 1, 4, and 5 hold. Additionally, suppose $g(x, \xi) = \max[h(x, \xi)]$ satisfies the above conditions and ϕ_k is a scalar smoothing function. Then $\hat{p}_k(\cdot)$ is \bar{L}_k -weakly convex for some constant \bar{L}_k that depends on $L'_h(\xi)$, $L'_{\phi, k}$, $M'_{\phi, k}$, and the diameter of X_ν .

Proof. First, note that $g(\cdot, \xi) := \max[h(\cdot, \xi)]$ is $L'_h(\xi)$ -weakly convex on X_ν for each $\xi \in \Xi$, see Lemma 4.2 of Drusvyatskiy and Paquette [22]. In particular, this implies that for any $y, z \in X_\nu$,

$$g(z, \xi) = \max[h(z, \xi)] \geq g(y, \xi) + s_y^T(\xi)(z - y) - \frac{L'_h(\xi)}{2} \|z - y\|^2$$

for any $s_y(\xi) \in \partial \max[h(y, \xi)]$. The Lipschitz continuity of $d\phi_k(\cdot)$ implies that for any scalars v and w :

$$\phi_k(w) \geq \phi_k(v) + d\phi_k(v)(w - v) - \frac{L'_{\phi, k}}{2} (w - v)^2.$$

From the monotonicity Assumption (4b) and by recursively applying the above result, we have for any

$\xi \in \Xi$ and $y, z \in X_\nu$:

$$\begin{aligned}
\phi_k(g(z, \xi)) &\geq \phi_k\left(g(y, \xi) + s_y^T(\xi)(z - y) - \frac{L'_h(\xi)}{2}\|z - y\|^2\right) \\
&\geq \phi_k\left(g(y, \xi) + s_y^T(\xi)(z - y)\right) - \frac{L'_h(\xi)}{2}d\phi_k\left(g(y, \xi) + s_y^T(\xi)(z - y)\right)\|z - y\|^2 - \frac{L'_{\phi,k}\left(L'_h(\xi)\right)^2}{8}\|z - y\|^4 \\
&\geq \phi_k(g(y, \xi)) + (d\phi_k(g(y, \xi)) s_y(\xi))^T(z - y) - \frac{L'_{\phi,k}}{2}\left(s_y^T(\xi)(z - y)\right)^2 - \frac{L'_h(\xi)M'_{\phi,k}}{2}\|z - y\|^2 - \\
&\quad \frac{L'_{\phi,k}\left(L'_h(\xi)\right)^2}{8}\|z - y\|^4 \\
&\geq \phi_k(g(y, \xi)) + (d\phi_k(g(y, \xi)) s_y(\xi))^T(z - y) - \frac{\|z - y\|^2}{2}\left[L'_{\phi,k}\|s_y(\xi)\|^2 + L'_h(\xi)M'_{\phi,k} + \right. \\
&\quad \left. \frac{L'_{\phi,k}\left(L'_h(\xi)\right)^2 \text{diam}(X_\nu)^2}{4}\right].
\end{aligned}$$

Taking expectation on both sides and noting that $\|s_y(\xi)\| \leq L'_h(\xi)$, we get

$$\begin{aligned}
\mathbb{E}[\phi_k(g(z, \xi))] &\geq \mathbb{E}[\phi_k(g(y, \xi))] + \mathbb{E}\left[(d\phi_k(g(y, \xi)) s_y(\xi))^T(z - y)\right] - \\
&\quad \frac{\|z - y\|^2}{2}\left[L'_{\phi,k}\mathbb{E}\left[\left(L'_h(\xi)\right)^2\right] + \mathbb{E}\left[L'_h(\xi)\right]M'_{\phi,k} + \frac{L'_{\phi,k}\mathbb{E}\left[\left(L'_h(\xi)\right)^2\right]\text{diam}(X_\nu)^2}{4}\right].
\end{aligned}$$

Therefore, \hat{p}_k is \bar{L}_k -weakly convex on X_ν with

$$\bar{L}_k := L'_{\phi,k}\mathbb{E}\left[\left(L'_h(\xi)\right)^2\right] + \mathbb{E}\left[L'_h(\xi)\right]M'_{\phi,k} + \frac{1}{4}L'_{\phi,k}\mathbb{E}\left[\left(L'_h(\xi)\right)^2\right]\text{diam}(X_\nu)^2. \quad \square$$

We now outline our proposal for estimating the weak convexity parameter \bar{L}_k of \hat{p}_k for the above setting. First, we note that estimate of the constants $L'_{\phi,k}$ and $M'_{\phi,k}$ can be obtained from Proposition 7. Next, we propose to replace the unknown constant $\text{diam}(X_\nu)$ with the diameter $2r$ of the sampling ball in Algorithm 4. Finally, we note that estimating the Lipschitz constant $L'_h(\xi)$ for any $\xi \in \Xi$ is tricky since it would involve looking at all of the $|J|$ constraints in general, just the thing we wanted to avoid! To circumvent this, we propose to replace $L'_h(\xi)$ in \bar{L}_k with an estimate of the local Lipschitz constant of our choice of the B -subdifferential element of $g(\cdot, \xi)$ through sampling (similar to the proposal in Algorithm 4). When g is defined through the solution of the recourse formulation considered, we can estimate a B -subdifferential element by appealing to linear programming duality and use this as a proxy for the ‘gradient’ of g . Note that a crude estimate of the step length does not affect the convergence guarantee of the stochastic subgradient method of Davis and Drusvyatskiy [19].

D Additional computational results

We present results of our replication experiments for the three case studies in the paper, and also consider a variant of Case study 2 that has a known analytical solution to benchmark our proposed approach.

Figure 5 presents the enclosure of the trajectories of the EF generated by ten replicates of the proposed approach when applied to Case study 1. For each value of the objective bound ν , we plot the smallest and largest risk level determined by Algorithm 2 at that bound over the different replicates. The reader can verify that the risk levels returned by the proposed algorithm do not vary significantly across the different replicates, with the maximum difference in the risk levels across the 26 points on the EF being a factor of 1.42. Figure 6 presents the corresponding plot for Case study 2. The reader can again verify that the risk levels returned by the proposed algorithm do not vary significantly across the different replicates, with the maximum difference in the risk levels across the 31 points on the EF being a factor of 1.21. Figure 7 presents the corresponding plot for Case study 3. The maximum difference in the risk levels at the 26 points on the EF for this case is a factor of 1.25 across the ten replicates.

Figure 5: Enclosure of the trajectories of the efficient frontier for Case study 1 generated by ten replicates of the proposed approach.

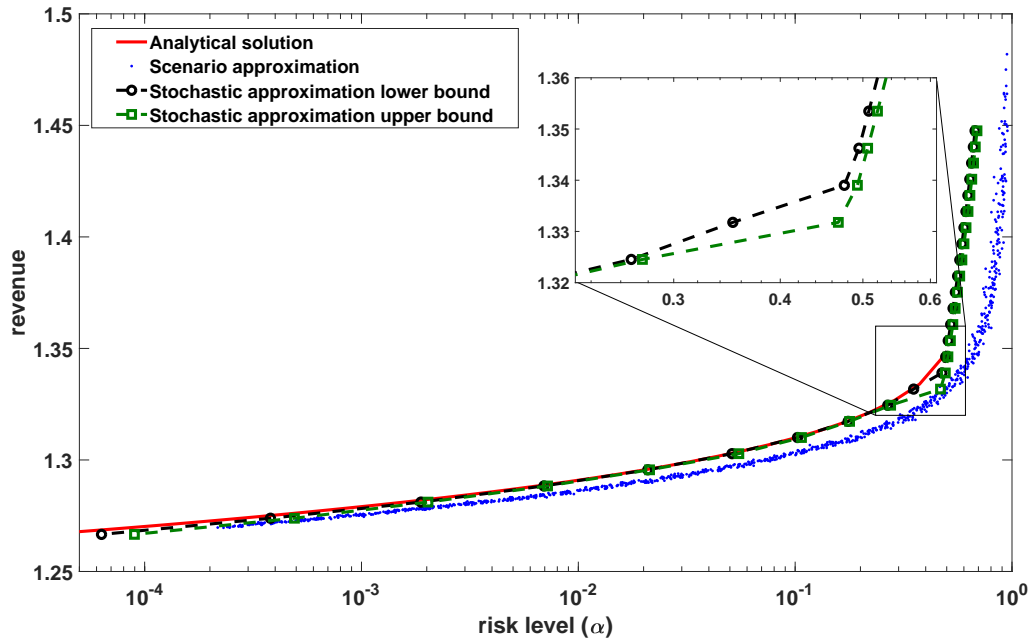


Figure 6: Enclosure of the trajectories of the efficient frontier for Case study 2 generated by ten replicates of the proposed approach.

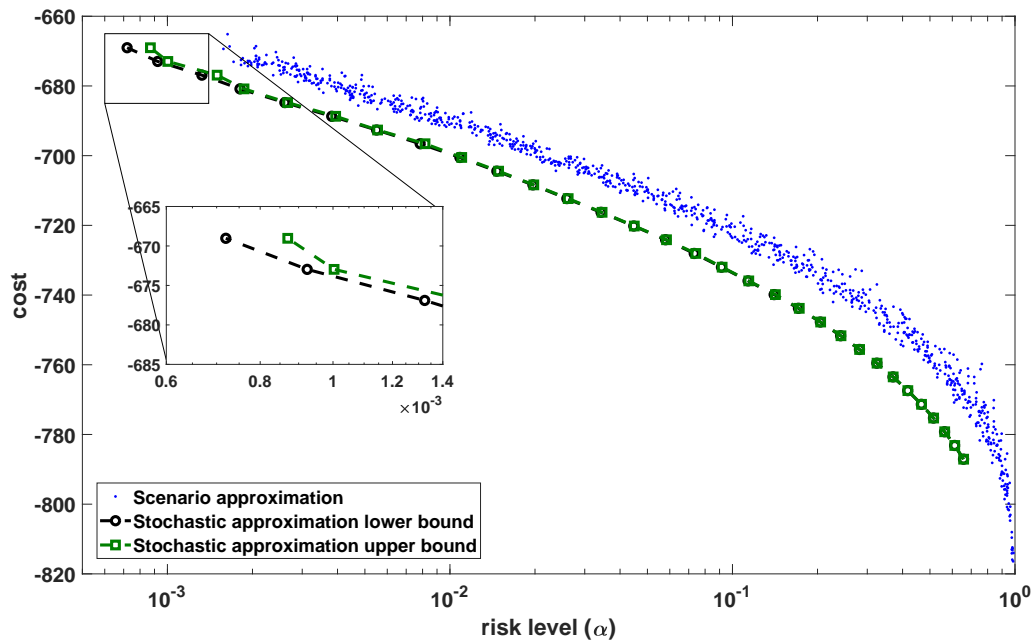
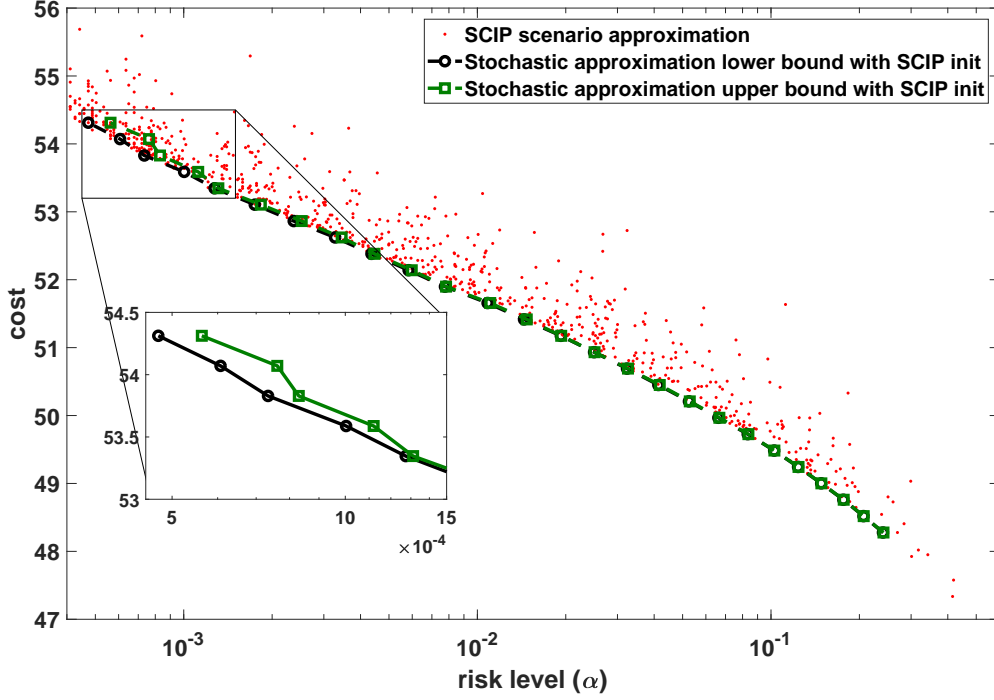


Figure 7: Enclosure of the trajectories of the efficient frontier for Case study 3 generated by ten replicates of the proposed approach.



Case study 4. We consider Case study 2 when the random variables $\xi_{ij} \sim \mathcal{P} := \mathcal{N}(0, 1)$ are i.i.d., the number of variables $n = 100$, the number of constraints $m = 100$, and bound $U = 100$. The EF can be computed analytically in this case, see Section 5.1.1 of Hong et al. [32]. Figure 8 compares a typical EF obtained using our approach against the analytical EF and the solutions generated by the tuned scenario approximation algorithm. Our proposed approach is able to converge to the analytical EF, whereas the scenario approximation method is only able to determine a suboptimal EF. Our proposal takes 1974 seconds on average (and a maximum of 2207 seconds) to approximate the EF using 32 points, whereas it took the scenario approximation a total of 16007 seconds to generate its 1000 points in Figure 8. We note that about 60% of the reported times for our method is spent in generating random numbers because the random variable ξ is high-dimensional.

Figure 9 presents an enclosure of the trajectories of the EF generated by the proposed approach over ten replicates for this case. The reader can verify that the risk levels returned by the proposed algorithm do not vary significantly across the different replicates, with the maximum difference in the risk levels across the 32 points on the EF being a factor of 1.002.

Figure 8: Comparison of the efficient frontiers for Case study 4.

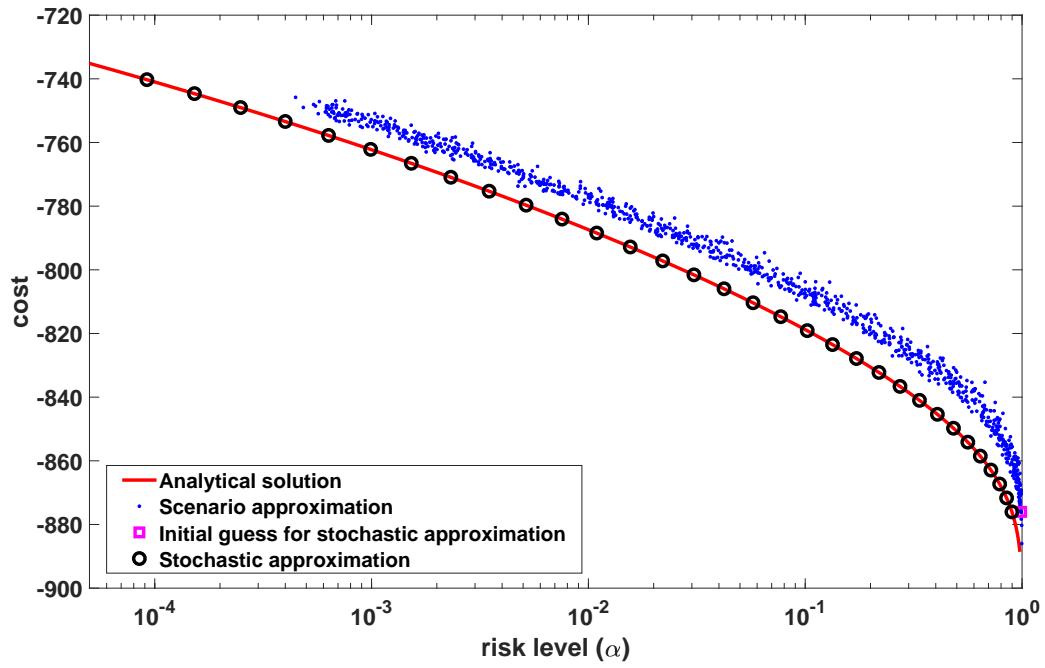


Figure 9: Enclosure of the trajectories of the efficient frontier for Case study 4 generated by ten replicates of the proposed approach.

