

1 **A SMOOTHING DIRECT SEARCH METHOD FOR MONTE CARLO-BASED BOUND**
2 **CONSTRAINED COMPOSITE NONSMOOTH OPTIMIZATION**

3 XIAOJUN CHEN*, C. T. KELLEY†, FENGMIN XU‡, AND ZAIKUN ZHANG§

January 31, 2018

4 **Abstract.** We propose and analyze a smoothing direct search algorithm for finding a minimizer of a nonsmooth nonconvex
5 function over a box constraint set, where the objective function values cannot be computed directly, but are approximated
6 by Monte Carlo simulation. In the algorithm, we adjust the stencil size, the sample size, and the smoothing parameter
7 simultaneously so that the stencil size goes to zero faster than the smoothing parameter and the square root of the sample size
8 goes to infinity faster than the reciprocal of the stencil size. We prove that with probability one any accumulation point of the
9 sequence generated by the algorithm is a Clarke stationary point. We report on numerical results from statistics and financial
10 applications.

11 **Key words.** Sampling methods, direct search algorithm, Monte Carlo simulation, non-smooth optimization, smoothing
12 functions, Clarke stationarity

13 **AMS subject classifications.** 65K05, 65K10, 90C30

14 **1. Introduction.** Let $G : R^{n+m} \rightarrow R$ be a locally Lipschitz continuous function, $F : R^n \rightarrow R^m$ be a
15 continuously differentiable function and $f : R^n \rightarrow R$ be a composite nonsmooth function of the form

16 (1.1)
$$f(x) = G(x, F(x)).$$

17 In this paper, we consider the nonsmooth minimization problem

18 (1.2)
$$\begin{aligned} \min \quad & f(x) \\ \text{s.t.} \quad & x \in X := \{x \in R^n \mid \ell \leq x \leq u\}, \end{aligned}$$

19 where $\ell, u \in R^n$, $-\infty < \ell_i < u_i < \infty$, $i = 1, \dots, n$ and the function value $F(x)$ is not computed exactly,
20 but rather approximated by a Monte Carlo simulation F^N , where N is the sample size of the Monte Carlo
21 simulation.

22 We assume that the function G admits a smoothing approximation in the sense of [10]. By this we
23 mean that there is a family of differentiable functions $\hat{G}(\cdot, \mu)$ which converges to G as $\mu \rightarrow 0$. We make
24 this precise in Definition 2.3. In this paper we use the smoothing parameter μ to design an algorithm which
25 takes gradients of \hat{G} to determine descent directions and then decreases μ as the iteration progresses.

26 Our use of the term nonsmooth in this context is standard [14] and refers to functions which, while not
27 differentiable (smooth), are locally Lipschitz continuous and the generalized derivatives are well-defined in
28 the sense of [14].

29 In [11] we considered a similar problem in the more general situation where the objective function was
30 not everywhere defined and capturing the domain of f was part of the problem. In this paper the objective
31 function is everywhere defined and can be approximated by a smoothing approach. The results in this paper
32 exploit the structure of that special case to simplify the analysis and improve the efficiency of the method
33 via smoothing.

34 A large class of nonsmooth functions have the form (1.1). The box constraint does not restrict appli-
35 cations where the objective functions have minimizers in a compact set. In other words, if the level set

*Department of Applied Mathematics, The Hong Kong Polytechnic University, Hung Hom, Kowloon, Hong Kong, China (maxjchen@polyu.edu.hk). The work of this author was partially supported by Hong Kong Research Grant Council grant PolyU153000/15p.

†North Carolina State University, Department of Mathematics, Box 8205, Raleigh, NC 27695-8205, USA (Tim.Kelley@ncsu.edu). The work of this author was partially supported by the Consortium for Advanced Simulation of Light Water Reactors (www.casl.gov), and Simulation of Nuclear Reactors under U.S. Department of Energy Contract No. DE-AC05-00OR22725, Army Research Office grant #W911NF-16-1-0504 and National Science Foundation Grant DMS-1406349.

‡School of Economics and Finance, Xi'an Jiaotong University, Xi'an, 710049, China (fengminxu@mail.xjtu.edu.cn). The work of this author was partially supported by the Chinese Natural Science Foundation 11571271, 71331001.

§Department of Applied Mathematics, The Hong Kong Polytechnic University, Hung Hom, Kowloon, Hong Kong, China (zaikun.zhang@polyu.edu.hk). The work of this author was supported by the start-up grant 1-ZVHT from The Hong Kong Polytechnic University and the Early Career Scheme grant PolyU 253012/17P from RGC Hong Kong.

36 $\{x : f(x) \leq f(x_0)\}$ of f with a point $x_0 \in R^n$ is bounded, then minimizing f over R^n can be equivalently
 37 written as a box constrained optimization problem (1.2). One example (for which $m = n$) is the expected
 38 value version of the stochastic variational inequality problem [12, 49]: Given the induced probability space
 39 $(\Xi \subset R^\ell, \mathcal{A}, \mathcal{P})$ and a convex set $\Omega \subseteq R^n$, find $x^* \in \Omega$ such that

$$40 \quad (1.3) \quad (x - x^*)^T F(x^*) \geq 0, \quad \forall x \in \Omega,$$

41 where $F(x) := E[\phi(\xi, x)]$, and $\phi : \Xi \times R^n \rightarrow R^n$ is continuously differentiable with respect to x for almost
 42 all $\xi \in \Xi$ and \mathcal{A} -measurable with respect to ξ . The stochastic variational inequality problem (1.3) reduces
 43 to the stochastic complementarity problem:

$$44 \quad (1.4) \quad x \geq 0, \quad F(x) \geq 0, \quad x^T F(x) = 0,$$

when $\Omega = R_+^n = \{x \in R^n \mid x \geq 0\}$, and the system of stochastic nonsmooth equations:

$$F(x) = 0$$

45 when $\Omega = R^n$. In this case, the approximation is via Monte Carlo simulation

$$46 \quad F(x) = E[\phi(\xi, x)] \approx F^N(x) := \frac{1}{N} \sum_{i=1}^N \phi(\xi^i, x),$$

47 where N is the sample size and $\xi^i, i = 1, \dots, N$ are observations of $\xi \in \Xi$.

48 We can express problem (1.3) as a minimization problem [20]

$$49 \quad (1.5) \quad \min_{x \in R^n} \|x - \text{Proj}_\Omega(x - F(x))\|_2^2$$

50 where Proj_Ω is the projection onto the set Ω . In this formulation the optimal function value is zero.

51 Another example of problem (1.2) is the ℓ_1 -norm regularized minimization problem

$$52 \quad (1.6) \quad \min_{x \in R^n} \|F(x)\|_1 + \lambda \|x\|_1,$$

53 where $F(x)$ is approximated by a Monte Carlo simulation.

54 If problems (1.5) and (1.6) have minimizers in a compact set, then there is a box constraint set X , which
 55 might be difficult to determine a priori, such that problems (1.5) and (1.6) can be equivalently written as
 56 problem (1.2) with $f(x) = \|x - \text{Proj}_\Omega(x - F(x))\|_2^2$ and $f(x) = \|F(x)\|_1 + \lambda \|x\|_1$, respectively.

57 We will exploit the composition structure by using a smoothing function for f . We will show that if f
 58 is replaced by the outcome of Monte Carlo simulation and one has full knowledge of the nonsmoothness, we
 59 can develop a smoothing direct search method with Monte Carlo simulation, which has global convergence
 60 to a Clarke stationary point of problem (1.2) with probability one (w.p.1.).

61 For example, if f has the form (1.1), we can define the smoothing function for f as

$$62 \quad (1.7) \quad \hat{f}(x, \mu) = \hat{G}(x, F(x), \mu)$$

63 and, when $F(x)$ is replaced by the Monte Carlo outcome $F^N(x)$, we set the smoothing Monte Carlo simulation
 64 as

$$65 \quad (1.8) \quad \tilde{f}(x, \mu, N) = \hat{G}(x, F^N(x), \mu).$$

66 We consider stencil-based direct search methods in this paper. By this we mean that at each step in the
 67 optimization the function is evaluated at a set of points of the form

$$68 \quad \{x \pm hv_i\}_{i=1}^n \cup \{x\}$$

69 where x is the current point and h is the stencil size. The directions v can be quite general [3, 28, 31, 32]. In
 70 this paper we will use the positive and negative coordinate directions, which is sufficient for our application.

71 For any fixed smoothing parameter $\mu > 0$, the function \hat{f} is continuously differentiable with respect to
 72 x . The main contribution of this paper is to propose a smoothing direct search algorithm with Monte Carlo
 73 simulation for solving problem (1.2) and prove the convergence of the algorithm when the stencil size h and
 74 smoothing parameter μ go to zero with the rate $h/\mu \rightarrow 0$, and the sample size N goes to infinity with the
 75 rate $(h\sqrt{N})^{-1} \rightarrow 0$.

76 Convergence analysis of direct search algorithms for smooth optimization problems where function values
 77 can be computed exactly have been well studied in [15, 16, 23, 28, 47]. Nonsmooth problems have been
 78 considered in [2–4, 11, 28]. Theory and algorithms for problems where the function evaluations require
 79 embedded Monte Carlo simulations have been carefully considered for optimization problems [11, 29, 33, 45,
 80 46, 49] and for nonlinear equations [55, 58]. The new algorithm in this paper exploits the structure of the
 81 problem and properties of smoothing methods to allow for using coordinate basis as fixed stencil search
 82 directions, simplifying the approaches of [3, 11, 28] for nonsmooth problems while preserving the convergence
 83 results.

84 Direct search methods have been coupled with randomized methods in [52] where the randomization was
 85 in the sampling and the optimization problem itself was deterministic. In [48] a generalized pattern search
 86 algorithm was applied to a problem where the objective function f was an expectation. The objective was a
 87 function of continuous and categorical variables and was assumed to be a smooth function of the continuous
 88 variables. Neither of these papers consider nonsmooth problems.

89 This paper is organized as follows. In section 2, we present a smoothing direct search algorithm for
 90 problem (1.2) where the function values $f(x)$ can be computed directly, and prove the convergence of the
 91 algorithm. In section 3, we extend the algorithm and convergence analysis to a smoothing direct search
 92 algorithm for (1.2) where the function values $f(x)$ cannot be computed directly, but are approximated
 93 by Monte Carlo simulation. In section 4, we present numerical experiments which include examples from
 94 statistical learning, and portfolio selection using test problems from the OR-Library [5] and real data from
 95 the Shanghai-Shenzhen stock market.

96 **2. A smoothing direct search method.** We begin by reviewing sampling direct search methods in
 97 the context of the smooth optimization problem. Let the set of search directions be an orthonormal basis
 98 $V = \{v_1, v_2, \dots, v_n\}$. Let h be the stencil size along those search directions. A stencil centered at x with h
 99 is the set of points $\{x \pm hv_i\}_{i=1}^n \cup \{x\}$. More general stencils can be used [3, 28, 31, 32] but are not needed
 100 for the applications in this paper.

101 The concept of stencil failure is important in both the algorithms and the analysis.

102 DEFINITION 2.1. *We say that stencil failure has occurred if*

$$103 \quad (2.1) \quad f(x) \leq f(x \pm hv_i) \quad \text{for } i = 1, \dots, n.$$

104 For simplicity, in this paper we will use

$$105 \quad (2.2) \quad V = \{e_1, e_2, \dots, e_n\},$$

106 for each iteration. Here e_i is the i -th coordinate vector. The algorithms and convergence analysis can be
 107 extended to an orthonormal basis.

108 It is easy to show [16, 27, 28] that if f is Lipschitz continuously differentiable in X , then (2.1) implies
 109 that

$$110 \quad (2.3) \quad \|\nabla f(x)\| = O(h)$$

111 uniformly for $x \in X$. To see this note that Lipschitz continuous differentiability of f in X and (2.1) imply
 112 that

$$113 \quad \frac{\partial f(x)}{\partial x_i} h + O(h^2) = f(x + he_i) - f(x) \geq 0, \quad i = 1, \dots, n$$

114 and

$$115 \quad -\frac{\partial f(x)}{\partial x_i} h + O(h^2) = f(x - he_i) - f(x) \geq 0, \quad i = 1, \dots, n$$

116 uniformly in X . Hence $\|\nabla f(x)\| = O(h)$, uniformly in X .

117 Sampling methods evaluate the objective function at the points of the stencil. If the current point is the
 118 best (stencil failure at the current point), then the stencil size is reduced. If the current point is not the best
 119 on the stencil, then the new best point becomes the current point. Algorithm `direct_search` is a version of
 120 the method for minimizing a continuously differentiable objective function f within a convex set X .

Algorithm `direct_search` (x, f, h)

```

for forever do
   $f_{base} = f(x)$ 
   $f_{min} = \min\{f(y) \mid y = x \pm hv, v \in V \text{ and } y \in X\}$ 
   $\hat{y} \in \{y \mid f(y) = f_{min}, y = x \pm hv, v \in V \text{ and } y \in X\}$ 
  if  $f_{min} \geq f_{base}$  then
     $h \leftarrow h/2$ 
  else
     $x \leftarrow \hat{y}$ 
  end if
end for

```

121 Algorithm `direct_search` generates two sequences $\{x_k\}$ and $\{h_k\}$. The convergence analysis is based
 122 on their subsequences $\{\tilde{x}_t\}$ and $\{\tilde{h}_t\}$, whose generation is described in the following box with initial points
 123 x_0 and h_0 , and $t = 0$.

Sequences generated by Algorithm `direct_search`

```

for  $k \geq 0$ 
   $\hat{y} \in \operatorname{argmin}\{f(y) \mid y = x_k \pm h_k v, v \in V \text{ and } y \in X\}$ 
  if  $f(\hat{y}) \geq f(x_k)$  (stencil failure)
    set  $x_{k+1} = x_k, h_{k+1} = h_k/2, t = t + 1, \tilde{x}_t = x_{k+1}, \tilde{h}_t = h_{k+1}$ 
  else
     $x_{k+1} = \hat{y}, h_{k+1} = h_k$ 
  end if
end for

```

125 In Algorithm `direct_search`, we must choose an initial point $x \in X$. This requirement for a feasible
 126 starting point in a box constraint set X is easy to satisfy. The worst case cost of a sweep through the stencils
 127 for a fixed $h > 0$ is sampling every point on the finite set

$$128 \quad \Omega_h(x) = \{x_0 + mhv \mid m = 1, 2, \dots \text{ and } v \in V\} \cup X,$$

129 where x_0 is either the initial point or the first point after h has been reduced. This worst case is, in our
 130 experience, very unlikely.

131 In our formulation the search is non-opportunistic. By this we mean that the minimization is done over
 132 the entire stencil. The analysis is the same for the opportunistic version, where the first point with a smaller
 133 function value than f_{base} is used. The reason is that the stencil size is only reduced when the stencil fails.
 134 Stencil failure can only take place if the entire stencil is sampled. Before then, it does not matter if the
 135 search is opportunistic or not.

136 The convergence proof of Algorithm `direct_search` is based on the stencil directions such that if stencil
 137 failure happens at the current point, then some type of approximate necessary condition holds. This idea can
 138 be made very general with different stencils and different smoothness requirements on the objective function
 139 f [2–4, 11, 28].

140 We consider the following first-order stationarity measure

$$141 \quad (2.4) \quad \chi(x) = \max_{x+d \in X, \|d\| \leq 1} [-d^T \nabla f(x)].$$

142 It is easy to check that, if $x \in X$ is a local minimizer of (1.2), then $\chi(x) = 0$.

143 PROPOSITION 2.2. Assume that f is Lipschitz continuously differentiable. Let $\{x_k\}$ with $x_0 \in X$ be the
 144 sequence generated by Algorithm `direct_search`. Then

$$145 \quad (2.5) \quad \liminf_{k \rightarrow \infty} \chi(x_k) = 0.$$

146 Moreover, stencil failure happens at infinitely many iterates, and for each limit point x of the stencil failure
 147 iterates, it holds that

$$148 \quad (2.6) \quad \chi(x) = 0.$$

149 *Proof.* Let h_0 be the initial stencil size and

$$150 \quad (2.7) \quad X_t = X \cap \left\{ x_0 + \sum_{i=1}^n \frac{j_i h_0}{2^t} e_i \mid j_i = 0, \pm 1, \pm 2, \dots \right\}, \quad t = 0, 1, 2, \dots$$

151 Since X is bounded, X_t is a finite set, and it contains at least one iterate. Stencil failure occurs at \tilde{x}_t , the
 152 last iterate contained in X_t and the size of the stencil at that iteration is $\tilde{h}_t = h_0/2^t$. For each $t \geq 1$, define

$$153 \quad (2.8) \quad I_t^+ = \{i \mid \tilde{x}_t + \tilde{h}_t e_i \in X, 1 \leq i \leq n\},$$

$$154 \quad (2.9) \quad I_t^- = \{i \mid \tilde{x}_t - \tilde{h}_t e_i \in X, 1 \leq i \leq n\},$$

156 and denote $g_t = \nabla f(\tilde{x}_t)$. Let L be the Lipschitz constant of ∇f . Then by the definition of stencil failure
 157 and Taylor's theorem,

$$158 \quad (2.10) \quad 0 \leq f(\tilde{x}_t + \tilde{h}_t e_i) - f(\tilde{x}_t) \leq \tilde{h}_t e_i^T g_t + \frac{L}{2} \tilde{h}_t^2 \quad \text{for all } i \in I_t^+,$$

$$159 \quad (2.11) \quad 0 \leq f(\tilde{x}_t - \tilde{h}_t e_i) - f(\tilde{x}_t) \leq -\tilde{h}_t e_i^T g_t + \frac{L}{2} \tilde{h}_t^2 \quad \text{for all } i \in I_t^-,$$

161 and consequently,

$$162 \quad (2.12) \quad e_i^T g_t \geq -\frac{L\tilde{h}_t}{2} \quad \text{for all } i \in I_t^+,$$

$$163 \quad (2.13) \quad e_i^T g_t \leq \frac{L\tilde{h}_t}{2} \quad \text{for all } i \in I_t^-.$$

165 By the assumption, there exists a positive constant Υ such that

$$166 \quad (2.14) \quad \|\nabla f(x)\| \leq \Upsilon \quad \text{for all } x \in X.$$

167 Specifically, $\|g_t\| \leq \Upsilon$. Therefore, for each d such that $\tilde{x}_t + d \in X$ and $\|d\| \leq 1$, it holds that

$$168 \quad (2.15) \quad \begin{aligned} -d^T g_t &= -\sum_{i=1}^n d_i (g_t)_i \\ &= -\sum_{i \in I_t^+ \setminus I_t^-} d_i (g_t)_i - \sum_{i \in I_t^- \setminus I_t^+} d_i (g_t)_i - \sum_{i \in I_t^+ \cap I_t^-} d_i (g_t)_i - \sum_{i \notin I_t^+ \cup I_t^-} d_i (g_t)_i \\ &\leq n \max \left\{ \frac{L\tilde{h}_t}{2}, \Upsilon \tilde{h}_t \right\} + n \frac{L\tilde{h}_t}{2} + n \Upsilon \tilde{h}_t \\ &\leq 3n \max \left\{ \frac{L\tilde{h}_t}{2}, \Upsilon \tilde{h}_t \right\}, \end{aligned}$$

169 where the first inequality uses the fact that for each d with $\tilde{x}_t + d \in X$, $i \notin I_t^-$ implies $d_i > \tilde{h}_t$ and $i \notin I_t^+$
 170 implies $d_i < \tilde{h}_t$, since X is a bounded box.

171 Hence, we obtain

$$172 \quad (2.16) \quad \chi(\tilde{x}_t) \leq 3n \max \left\{ \frac{L\tilde{h}_t}{2}, \Upsilon\tilde{h}_t \right\} \rightarrow 0 \quad \text{when } t \rightarrow \infty.$$

173 Since $\{\tilde{x}_t\}$ is a subsequence of $\{x_k\}$ and X is bounded, we conclude that

$$174 \quad (2.17) \quad \liminf_{k \rightarrow \infty} \chi(x_k) = 0.$$

175 If x is an accumulation point of the stencil failure iterates $\{\tilde{x}_t\}$, then continuity of χ implies that $\chi(x) = 0$. \square

176 **2.1. Nonsmooth f .** In this subsection, we consider problem (1.1) where F can be computed exactly.
 177 We will use smoothing methods which approximate f by a parameterized family of smoothing functions
 178 $\hat{f}(\cdot, \mu)$ given by (1.7), where $\mu > 0$ is the smoothing parameter.

179 We formally give a definition of smoothing functions used in this paper.

180 **DEFINITION 2.3.** [10] Let $f : R^n \rightarrow R$ be a locally Lipschitz continuous function. We call $\hat{f} : R^n \times$
 181 $(0, \infty) \rightarrow R$ a smoothing function of f , if $\hat{f}(\cdot, \mu)$ is continuously differentiable and $\nabla \hat{f}(\cdot, \mu)$ is Lipschitz
 182 continuous in R^n for any fixed $\mu \in (0, \infty)$, and

$$183 \quad (2.18) \quad \lim_{x \rightarrow \hat{x}, \mu \downarrow 0} \hat{f}(x, \mu) = f(\hat{x}).$$

184 The limit in (2.18) is simultaneous in x and μ for all sequences $x_k \rightarrow \hat{x}$ and $\mu_k \rightarrow 0$ ($\mu_k \geq 0$).

185 Throughout this subsection we let $\nabla \hat{f}$ denote the gradient \hat{f} with respect to x .

186 **ASSUMPTION 2.1.** (i) There are constants $c_1, c_2 \geq 0$ such that for any $x \in R^n$, $\mu \in (0, 1]$,

$$187 \quad (2.19) \quad |f(x) - \hat{f}(x, \mu)| \leq \mu(c_1 + c_2|f(x)|).$$

188 (ii) \hat{f} satisfies the gradient consistency condition,

$$189 \quad (2.20) \quad \partial f(x) = \text{con}\{v \mid \nabla \hat{f}(x_k, \mu_k) \rightarrow v, \text{ for } x_k \rightarrow x, \mu_k \downarrow 0\},$$

190 where “con” denotes the convex hull and $\partial f(x)$ is the Clarke subgradient at x .

191 (iii) There are $\Upsilon > 0$, $\Gamma > 0$ and $\mu_- > 0$ such that $\|\nabla \hat{f}(x, \mu)\| \leq \Upsilon$ and

$$192 \quad (2.21) \quad \|\nabla \hat{f}(x, \mu) - \nabla \hat{f}(y, \mu)\| \leq \frac{\Gamma}{\mu} \|x - y\|$$

193 uniformly in $x, y \in X$, and $\mu \in (0, \mu_-)$.

In Assumption 2.1, $c_1, c_2, \Upsilon, \Gamma$ and μ_- are fixed constants which are independent of x . Since X is bounded and f is continuous, Assumption 2.1 (i) implies that there is a constant C such that

$$|f(x) - \hat{f}(x, \mu)| \leq \mu C, \quad \text{for } x \in X, \mu \in (0, 1]$$

194 which means that \hat{f} converges to f uniformly as $\mu \rightarrow 0$.

195 In section 4, we use examples to illustrate the definition of smoothing functions and Assumption 2.1.

196 Note that X is bounded and hence there are only finitely many points in the stencil for each h . Therefore
 197 the stencil will fail infinitely often.

198 In the case where f is known exactly and there is no embedded Monte Carlo simulation, we propose
 199 a smoothing direct search algorithm, Algorithm `smoothing_search` that decreases μ and h simultaneously,
 200 but in a way that ensures $h/\mu \rightarrow 0$ as $\mu \rightarrow 0$, which will be important in the convergence analysis.

201 Algorithm `smoothing_search` generates three sequences $\{x_k\}$, $\{h_k\}$ and $\{\mu_k\}$. The convergence analysis
 202 is based on their subsequences $\{\tilde{x}_t\}$, $\{\tilde{h}_t\}$ and $\{\tilde{\mu}_t\}$, whose generation is described in the following box
 203 with initial points x_0, h_0 , and μ_0 , and $t = 0$.

Algorithm `smoothing_search`(x, \hat{f}, h, μ, τ)

```

for forever do
   $\hat{f}_{base} = \hat{f}(x, \mu)$ 
   $\hat{f}_{min} = \min\{\hat{f}(y, \mu) \mid y = x \pm hv, v \in V \text{ and } y \in X\}$ 
   $\hat{y} \in \{y \mid \hat{f}(y, \mu) = \hat{f}_{min}, y = x \pm hv, v \in V \text{ and } y \in X\}$ 
  if  $\hat{f}_{min} \geq \hat{f}_{base}$  then
     $h \leftarrow h/2; \quad \mu \leftarrow \mu/2^\tau$ 
  else
     $x \leftarrow \hat{y}$ 
  end if
end for

```

Sequences generated by Algorithm `smoothing_search`

```

for  $k \geq 0$ 
   $\hat{y} \in \operatorname{argmin}\{\hat{f}(y, \mu) \mid y = x_k \pm h_k v, v \in V \text{ and } y \in X\}$ 
  if  $\hat{f}(\hat{y}, \mu_k) \geq \hat{f}(x_k, \mu_k)$  (Stencil failure)
    set  $x_{k+1} = x_k, h_{k+1} = h_k/2, \mu_{k+1} = \mu_k/2^\tau, t = t + 1, \tilde{x}_t = x_{k+1}, \tilde{h}_t = h_{k+1}, \tilde{\mu}_t = \mu_{k+1}$ 
  else
     $x_{k+1} = \hat{y}, h_{k+1} = h_k, \mu_{k+1} = \mu_k, k = k + 1$ 
  end if
end for

```

205 In Algorithm `smoothing_search`, $\tau \in (0, 1)$ is an input parameter. We must choose an initial point
 206 $x \in X$, the initial stencil size $h > 0$ and the initial smoothing parameter $\mu > 0$.

207 As an extension of (2.4), we use

$$208 \quad (2.22) \quad \tilde{\chi}(x) = \min_{v \in \partial f(x)} \left(\max_{x+d \in X, \|d\| \leq 1} -d^T v \right)$$

209 to measure the first-order stationarity of x with respect to problem (1.2) when f is locally Lipschitz continuous
 210 but not necessarily differentiable, where $\partial f(x)$ is the Clarke subdifferential of f at x [14, 41]. If f is smooth,
 211 then $\tilde{\chi}(\cdot)$ is the same as $\chi(\cdot)$. Moreover, if x is a local minimizer of problem (1.2), then there exists a
 212 $v \in \partial f(x)$ such that

$$213 \quad \max_{x+d \in X, \|d\| \leq 1} [-d^T v] = 0,$$

214 that is, $\tilde{\chi}(x) = 0$.

215 The convergence result follows the same argument as in the proof of Proposition 2.2 and Assumption
 216 2.1 on the smoothing function of f .

217 **THEOREM 2.4.** *Assume that Assumption 2.1 holds. Let $\{x_k, \mu_k\}$ with $x_0 \in X$ and $\mu_0 > 0$ be the iterates*
 218 *generated by Algorithm* `smoothing_search`, *and*

$$219 \quad (2.23) \quad \chi_k(x) = \max_{x+d \in X, \|d\| \leq 1} [-d^T \nabla \hat{f}(x, \mu_k)].$$

220 *Then*

$$221 \quad (2.24) \quad \liminf_{k \rightarrow \infty} \chi_k(x_k) = 0.$$

222 *Moreover, stencil failure happens at infinitely many iterates, and for each limit point x of the stencil failure*
 223 *iterates, it holds that*

$$224 \quad (2.25) \quad \tilde{\chi}(x) = 0.$$

225 *Proof.* Define X_t by (2.7). As in the proof of Proposition 2.2, we denote the last iterate in X_t , where
 226 stencil failur occurs by \tilde{x}_t , the corresponding stencil size by $\tilde{h}_t = h_0/2^t$, and the corresponding smoothing
 227 parameter by $\tilde{\mu}_t = \tilde{h}_t^2$. For each $t \geq 1$, define I_t^+ and I_t^- in the same way as in the proof of Proposition 2.2,
 228 and denote $\hat{g}_t = \nabla \hat{f}(\tilde{x}_t, \tilde{\mu}_t)$. According to the definition of stencil failure and Taylor expansion, noticing
 229 part (iii) of Assumption 2.1, we have

$$230 \quad (2.26) \quad 0 \leq \hat{f}(\tilde{x}_t + \tilde{h}_t e_i, \tilde{\mu}_t) - \hat{f}(\tilde{x}_t, \tilde{\mu}_t) \leq \tilde{h}_t e_i^T \hat{g}_t + \frac{\Gamma}{2\tilde{\mu}_t} \tilde{h}_t^2 \quad \text{for all } i \in I_t^+,$$

$$231 \quad (2.27) \quad 0 \leq \hat{f}(\tilde{x}_t - \tilde{h}_t e_i, \tilde{\mu}_t) - \hat{f}(\tilde{x}_t, \tilde{\mu}_t) \leq -\tilde{h}_t e_i^T \hat{g}_t + \frac{\Gamma}{2\tilde{\mu}_t} \tilde{h}_t^2 \quad \text{for all } i \in I_t^-,$$

233 and consequently,

$$234 \quad (2.28) \quad e_i^T \hat{g}_t \geq -\frac{\Gamma \tilde{h}_t}{2\tilde{\mu}_t} \quad \text{for all } i \in I_t^+,$$

$$235 \quad (2.29) \quad e_i^T \hat{g}_t \leq \frac{\Gamma \tilde{h}_t}{2\tilde{\mu}_t} \quad \text{for all } i \in I_t^-.$$

237 By Assumption 2.1, there exists a positive constant Υ such that $\|\hat{g}_t\| \leq \Upsilon$. Using similar argument to
 238 those for (2.15) and (2.16), we have

$$239 \quad (2.30) \quad \max_{\tilde{x}_t + d \in X, \|d\| \leq 1} -d^T \hat{g}_t \leq 3n \max \left\{ \frac{\Gamma \tilde{h}_t}{2\tilde{\mu}_t}, \Upsilon \tilde{h}_t \right\}.$$

240 Noticing the fact that $\tilde{h}_t/\tilde{\mu}_t \rightarrow 0$, we have

$$241 \quad (2.31) \quad \max_{\tilde{x}_t + d \in X, \|d\| \leq 1} -d^T \hat{g}_t \rightarrow 0 \quad \text{when } t \rightarrow \infty,$$

242 which implies (2.24).

243 Let x be a limit point of $\{\tilde{x}_t\}$, and $\{\tilde{x}_{t_i}\}$ be a subsequence that converges to x . Since $\{\hat{g}_t\}$ is bounded, we
 244 may suppose that $\{\hat{g}_{t_i}\}$ converges to a point v (if not, replace $\{t_i\}$ by an appropriately chosen subsequence).
 245 Let

$$246 \quad (2.32) \quad d^*(v) \in \operatorname{argmax}_{x+d \in X, \|d\| \leq 1} [-d^T v],$$

247 and

$$248 \quad (2.33) \quad y_i = \frac{x - \tilde{x}_{t_i} + d^*(v)}{\max\{\|x - \tilde{x}_{t_i} + d^*(v)\|, 1\}}.$$

249 Then $\|y_i\| \leq 1$, and $\tilde{x}_{t_i} + y_i \in X$ due to the convexity of X ($\tilde{x}_{t_i} + y_i$ lies on the line segment between \tilde{x}_{t_i}
 250 and $x + d^*$). Hence

$$251 \quad (2.34) \quad \begin{aligned} 0 &\leq \max_{x+d \in X, \|d\| \leq 1} -d^T v = -(d^*(v))^T v \\ &= \lim_{i \rightarrow \infty} -y_i^T \hat{g}_{t_i} \\ &\leq \lim_{i \rightarrow \infty} \max_{\tilde{x}_{t_i} + d \in X, \|d\| \leq 1} -d^T \hat{g}_{t_i} \\ &= 0. \end{aligned}$$

252 Notice that $v \in \partial f(x)$ according to the gradient consistency of \hat{f} . By the definition (2.22) of $\tilde{\chi}(\cdot)$, we have

$$253 \quad (2.35) \quad 0 \leq \tilde{\chi}(x) \leq \max_{x+d \in X, \|d\| \leq 1} -d^T v = 0,$$

254 which completes the proof. \square

255 **3. Smoothing direct search method with Monte Carlo simulations.** In this section we extend
 256 the algorithms and analysis from § 2 to the case where f is nonsmooth and approximated by a Monte Carlo
 257 simulation. Deterministic direct search methods for nonsmooth optimization problems have been studied
 258 in [2–4, 11, 23, 28, 42].

259 In this section, we assume that for any $x \in X$, $\mu > 0$, we can estimate the value of $\hat{f}(x, \mu)$ by Monte
 260 Carlo simulation $\tilde{f}(x, \mu, N)$ with N realizations. The value of $\tilde{f}(x, \mu, N)$ is random and the sample of N
 261 realizations is independently identically distributed (iid). We can view $\tilde{f}(x, \mu, N)$ as defined on a common
 262 probability space (see [49, pp. 156] for details).

263 The following is an assumption on the effectiveness of $\tilde{f}(\cdot, \cdot, N)$ as an approximation of $\hat{f}(\cdot, \cdot)$.

264 ASSUMPTION 3.1. For each $p \in (0, 1/2)$, there exist constants $\delta \in (0, 1)$, $c_F > 0$, $\bar{N} > 0$, and $\bar{\mu} > 0$ such
 265 that

$$266 \quad (3.1) \quad \text{Prob} \left(\sup_{x \in X} |\hat{f}(x, \mu) - \tilde{f}(x, \mu, N)| \geq \frac{c_f}{N^p} \right) \leq \delta$$

267 for each $N \geq \bar{N}$ and $\mu \in (0, \bar{\mu}]$.

268 Consider the composite nonsmooth function in the form (1.1) with

$$269 \quad (3.2) \quad F(x) = E[\phi(\xi, x)], \quad x \in X,$$

270 where ξ is a random vector. Let

$$271 \quad (3.3) \quad \tilde{F}^N(x) = \frac{1}{N} \sum_{i=1}^N \phi(\xi_i, x),$$

272 $\xi_1, \xi_2, \dots, \xi_N$ being iid samples of ξ . Assume that $\phi(\xi, x)$ is *sub-exponential* for each $x \in X$ (see Appendix B
 273 for the definition of sub-exponential random variables/vectors), and that $\phi(\xi, \cdot)$ is L -Lipschitz continuous
 274 with respect to $x \in X$ for a constant L independent of ξ . Then, as we show in Appendix B, for any
 275 $p \in (0, 1/2)$, there exist constants $\delta \in (0, 1)$, $c_F > 0$ and $\bar{N} > 0$ such that

$$276 \quad (3.4) \quad \text{Prob} \left(\sup_{x \in X} \|F(x) - \tilde{F}^N(x)\| \geq \frac{c_F}{N^p} \right) \leq \delta$$

for each $N \geq \bar{N}$. If G is Lipschitz continuous with respect to $F(x)$ and \hat{G} satisfies Assumption 2.1, then
 there exists a constant L_F independent of N and μ such that

$$|\hat{G}(x, F(x), \mu) - \hat{G}(x, \tilde{F}^N(x), \mu)| \leq L_F \|F(x) - \tilde{F}^N(x)\|.$$

This, together with (3.4), implies that function

$$\tilde{f}(x, \mu, N) \equiv \hat{G}(x, \tilde{F}^N(x), \mu)$$

277 fulfills Assumption 3.1 with $c_f = L_F c_F$.

278 Algorithm `mc_smoothing_search` for the embedded Monte Carlo case is a simple extension of Algorithm
 279 `smoothing_search`.

280 In Algorithm `mc_smoothing_search`, $\tau \in (0, 1)$ and $\gamma > 1$ are input parameters. The objective function
 281 f is evaluated through $\tilde{f}(x, \mu, N)$, the Monte Carlo simulation of $\hat{f}(x, \mu)$ with sample size N , where $\hat{f}(x, \mu)$
 282 is a smoothing function of f that is defined in Definition 2.3 and satisfies Assumption 2.1.

283 Algorithm `mc_smoothing_search` generates four sequences $\{x_k\}$, $\{h_k\}$, $\{\mu_k\}$ and $\{N_k\}$. The convergence
 284 analysis is based on their subsequences $\{\tilde{x}_t\}$, $\{\tilde{h}_t\}$, $\{\tilde{\mu}_t\}$ and $\{\tilde{N}_t\}$, whose generation is described in the
 285 following box with initial points x_0, h_0, μ_0 and N_0 , and $t = 0$.

Algorithm mc_smoothing_search ($x, \tilde{f}, h, \mu, N, \tau, \gamma$)

```

for forever do
   $\tilde{f}_{base} = \tilde{f}(x, \mu, N)$ 
   $\tilde{f}_{min} = \min\{\tilde{f}(y, \mu, N) \mid y = x \pm hv, v \in V \text{ and } y \in X\}$ 
   $\hat{y} \in \{y \mid \tilde{f}(y, \mu, N) = \tilde{f}_{min}, y = x \pm hv, v \in V \text{ and } y \in X\}$ 
  if  $\tilde{f}_{min} \geq \tilde{f}_{base}$  then
     $h \leftarrow h/2;$   $\mu \leftarrow \mu/2^\tau;$   $N \leftarrow 4^\gamma N$ 
  else
     $x \leftarrow \hat{y}$ 
  end if
end for

```

Sequences generated by Algorithm mc_smoothing_search
for $k \geq 0$
 $\hat{y} \in \operatorname{argmin}\{\tilde{f}(y, \mu, N) \mid y = x_k \pm h_k v, v \in V \text{ and } y \in X\}$
if $\tilde{f}(\hat{y}, \mu_k, N_k) \geq \tilde{f}(x_k, \mu_k, N_k)$ (**Stencil failure**)
 set $x_{k+1} = x_k, h_{k+1} = h_k/2, \mu_{k+1} = \mu_k/2^\tau, N_{k+1} = 4^\gamma N_k,$
 $t = t + 1, \tilde{x}_t = x_{k+1}, \tilde{h}_t = h_{k+1}, \tilde{\mu}_t = \mu_{k+1}, \tilde{N}_t = N_{k+1}$
else
 $x_{k+1} = \hat{y}, h_{k+1} = h_k, \mu_{k+1} = \mu_k, N_{k+1} = N_k, k = k + 1$
end if
end for

286

287 As in Proposition 2.2 and Theorem 2.4, the boundedness of X ensures the convergence of the algorithm.
 288 Let $x_0 \in X$ denote the initial point of Algorithm mc_smoothing_search, $h_0 > 0$ the initial stencil size, $\mu_0 > 0$
 289 the initial smoothing parameter, and $N_0 > 0$ the initial sample size.

290 The main result of the paper is Theorem 3.1, which states that w.p.1. the iteration has an accumulation
 291 point which is a Clarke stationary point.

292 **THEOREM 3.1.** *Suppose that the Monte Carlo simulations in Algorithm mc_smoothing_search are mu-*
 293 *tually independent for different N . Assume that Assumptions 2.1 and 3.1 hold. Let $\{x_k, \mu_k, N_k\}$ be the*
 294 *sequence generated by Algorithm mc_smoothing_search. Then*

295 (3.5)
$$\operatorname{Prob}\left(\liminf_{k \rightarrow \infty} \chi_k(x_k) = 0\right) = 1,$$

296 where χ_k is defined in (2.23), and

297 (3.6)
$$\operatorname{Prob}(\{x_k\} \text{ has an accumulation point } x \text{ such that } \tilde{\chi}(x) = 0) = 1.$$

298 *Proof.* Define X_t by (2.7). As before we denote the last iterate in X_t , where stencil failure occurs by \tilde{x}_t ,
 299 the corresponding stencil size by $\tilde{h}_t = h_0/2^t$, and the corresponding smoothing parameter by $\tilde{\mu}_t = \mu_0/2^{t\tau}$.
 300 The sample size at this point in the algorithm is $\tilde{N}_t = 4^{t\gamma} N_0$. Define index sets I_t^+ and I_t^- in the same way
 301 as in the proof of Proposition 2.2. Since stencil failure happens at \tilde{x}_t , we have

302 (3.7)
$$0 \leq \tilde{f}(\tilde{x}_t + \tilde{h}_t e_i, \tilde{\mu}_t, \tilde{N}_t) - \tilde{f}(\tilde{x}_t, \tilde{\mu}_t, \tilde{N}_t) \quad \text{for all } i \in I_t^+,$$

303 (3.8)
$$0 \leq \tilde{f}(\tilde{x}_t - \tilde{h}_t e_i, \tilde{\mu}_t, \tilde{N}_t) - \tilde{f}(\tilde{x}_t, \tilde{\mu}_t, \tilde{N}_t) \quad \text{for all } i \in I_t^-.$$

305

306 Let p be a constant such that

307 (3.9)
$$\frac{1}{2^\gamma} < p < \frac{1}{2}.$$

308 Then $\tilde{h}_t \tilde{N}_t^p \rightarrow \infty$ as $t \rightarrow \infty$. Set $\delta \in (0, 1)$ and $c_f > 0$ to be the constants that fulfill Assumption 3.1, and
 309 consider the event

$$310 \quad (3.10) \quad E_t = \left\{ \sup_{x \in X} |\hat{f}(x, \tilde{\mu}_t) - \tilde{f}(x, \tilde{\mu}_t, \tilde{N}_t)| \leq \frac{c_f}{\tilde{N}_t^p} \right\}.$$

311 By assumption, $\{E_t\}_{t=1}^\infty$ are mutually independent, and

$$312 \quad \text{Prob}(E_t) \geq 1 - \delta > 0$$

313 for each t sufficiently large so that $\tilde{N}_t \geq \bar{N}$ and $\tilde{\mu}_t \leq \bar{\mu}$. Therefore,

$$314 \quad (3.11) \quad \text{Prob}(E_t \text{ happens for infinitely many } t) = 1.$$

315 When E_t happens, according to (3.7) and (3.8), we have

$$316 \quad (3.12) \quad -\frac{2c_f}{\tilde{N}_t^p} \leq \hat{f}(\tilde{x}_t + \tilde{h}_t e_i, \tilde{\mu}_t) - \hat{f}(\tilde{x}_t, \tilde{\mu}_t) \leq \tilde{h}_t e_i^T \hat{g}_t + \frac{\Gamma}{2\tilde{\mu}_t} \tilde{h}_t^2 \quad \text{for all } i \in I_t^+,$$

$$317 \quad (3.13) \quad -\frac{2c_f}{\tilde{N}_t^p} \leq \hat{f}(\tilde{x}_t - \tilde{h}_t e_i, \tilde{\mu}_t) - \hat{f}(\tilde{x}_t, \tilde{\mu}_t) \leq -\tilde{h}_t e_i^T \hat{g}_t + \frac{\Gamma}{2\tilde{\mu}_t} \tilde{h}_t^2 \quad \text{for all } i \in I_t^-,$$

318 where $\hat{g}_t = \nabla \hat{f}(\tilde{x}_t, \tilde{\mu}_t)$ and Γ is the constant in item (iii) of Assumption 2.1. Subsequently, it holds that

$$320 \quad (3.14) \quad e_i^T \hat{g}_t \geq -\frac{\Gamma \tilde{h}_t}{2\tilde{\mu}_t} - \frac{2c_f}{\tilde{h}_t \tilde{N}_t^p} \quad \text{for all } i \in I_t^+,$$

$$321 \quad (3.15) \quad e_i^T \hat{g}_t \leq \frac{\Gamma \tilde{h}_t}{2\tilde{\mu}_t} + \frac{2c_f}{\tilde{h}_t \tilde{N}_t^p} \quad \text{for all } i \in I_t^-.$$

323 By Assumption 2.1, there exists a positive constant Υ such that $\|\hat{g}_k\| \leq \Upsilon$. Using similar argument as (2.15)
 324 and (2.16), we obtain from (3.14) and (3.15) that

$$325 \quad (3.16) \quad \max_{\tilde{x}_t + d \in X, \|d\| \leq 1} -d^T \hat{g}_t \leq 3n \max \left\{ \frac{\Gamma \tilde{h}_t}{2\tilde{\mu}_t} + \frac{2c_f}{\tilde{h}_t \tilde{N}_t^p}, \Upsilon \tilde{h}_t \right\}.$$

326 Hence, by (3.11) and the fact that $\tilde{h}_t \rightarrow 0$, $\tilde{h}_t/\tilde{\mu}_t \rightarrow 0$, and $\tilde{h}_t \tilde{N}_t^p \rightarrow \infty$, we have

$$327 \quad (3.17) \quad \text{Prob}(\{\chi_k(x_k)\} \text{ has a subsequence that converges to zero}) = 1,$$

328 which implies (3.5).

329 When $\liminf_{k \rightarrow \infty} \chi_k(x_k) = 0$, let $\{k_i\}$ be the index sequence such that $\chi_{k_i}(x_{k_i}) \rightarrow 0$. Since $\{x_{k_i}\}$ is
 330 bounded (guaranteed the boundedness of X), it has an accumulation point x . By the same argument that
 331 leads to the second part of Theorem 2.4, we have that $\tilde{\chi}(x) = 0$. Thus (3.6) holds. \square

332 **4. Numerical experiments.** In this section, we test Algorithm `mc_smoothing_search` on two prob-
 333 lems: a stochastic optimization problem arising from censored regression and a two-stage optimization
 334 problem arising from portfolio management. The problems in § 4.1 and 4.2.1 are derived from applications,
 335 but use synthetic data to enable us to control the sample size.

336 **4.1. Censored regression.** We consider the following regularized censored regression problem [1, 6,
 337 34, 35, 51, 53]

$$338 \quad (4.1) \quad \begin{aligned} & \min_{x \in R^n} f(x) \\ & \text{s.t. } -e \leq x \leq e, \end{aligned}$$

339 where $e \in R^n$ is the vector with all its entries being one and

$$340 \quad (4.2) \quad f(x) = E_{c,y}[(\max(c^T x, 0) - y)^2] + \lambda \sum_{i=1}^n \log(1 + |(x)_i|).$$

341 Here the random variable pair (c, y) represents a data set of interest ($c \in R^n$, $y \in R$), and $\lambda > 0$ is a
342 regularization parameter.

343 The regularization term

$$344 \quad \lambda \sum_{i=1}^n \log(1 + |(x)_i|)$$

345 in the objective function is used to enforce sparsity.

346 We assume that $c \sim N(0, I)$, and $y = \max(c^T x^* + \epsilon, 0)$ for some underlying ground truth feature x^* and
347 unobservable noise $\epsilon \sim N(0, \sigma^2)$. Moreover, we assume x^* is sparse, that is, x^* has few nonzero entries. Using
348 the concave regularized model (4.1), we want to recover a sparse feature x to approximate x^* as accurate as
349 possible given that $x^* \in \{x \mid \|x\|_\infty \leq 1\} = [-e, e]$.

The functions $(\max(c^T x, 0) - y)^2$ and $\log(1 + |(x)_i|)$ are not differentiable but do admit smoothing
functions (see Appendix A). Using the smoothing functions, we can define a smoothing function $\hat{f}_{c,y}(x, \mu)$
for $(\max(c^T x, 0) - y)^2$ which satisfies Assumption 2.1. From the convexity of $(\max(c^T x, 0) - y)^2$, the Clarke
subdifferential and the expectation can be exchanged, that is,

$$\partial E_{c,y}[(\max(c^T x, 0) - y)^2] = E_{c,y}[\partial(\max(c^T x, 0) - y)^2]$$

350 (see [14]). Moreover, $E_{c,y}[\hat{f}_{c,y}(x, \mu)]$ is a smoothing function for $E_{c,y}[(\max(c^T x, 0) - y)^2]$, and satisfies
351 Assumption 2.1 (see [8]). Problem (4.1) is a constrained nonsmooth nonconvex optimization problem where
352 the objective function values cannot be computed directly.

353 In practice the data in these problems are limited. To mimic the finite size of the data we will pose
354 an approximation to problem (4.1) that replaces the expectation with the sample average of a finite, but
355 large, data set. We will manage the sampling in the algorithm by randomly sampling from that data set.
356 To this end, we consider $X = \{x \mid \|x\|_\infty \leq 1\} = [-e, e]$, and we randomly generate a true feature $x^* \in R^{20}$
357 whose 5 nonzero entries are from uniform distribution on $[-1, 1]$. Independently, we generate samples c_i
358 from $c \sim N(0, I)$ and ϵ_i from $\epsilon \sim N(0, 0.01)$ with a sample size 10^7 . Let $y_i = \max(c_i^T x^* + \epsilon_i, 0)$. The new
359 problem is an approximation to (4.1) with a finite data set. We have

$$360 \quad (4.3) \quad \bar{f}(x, 10^7) = \frac{1}{10^7} \sum_{i=1}^{10^7} [(\max(c_i^T x, 0) - y_i)^2] + \lambda \sum_{i=1}^n \log(1 + |(x)_i|).$$

361 We used the regularization parameter $\lambda = 10^{-2}$. This is large enough to capture the sparsity exactly and
362 small enough to allow us to observe several iterations before the iteration stagnates.

363 We configure the optimization as follows:

- 364 • The algorithmic parameters are $c = 2$, $\gamma = 1.5$, and $\tau = 0.5$.
- 365 • We begin with $h = 0.5$ and terminate when $h \leq 10^{-3}$.
- 366 • $N = 100$ and $\mu = 0.1$ at the beginning of the iteration.

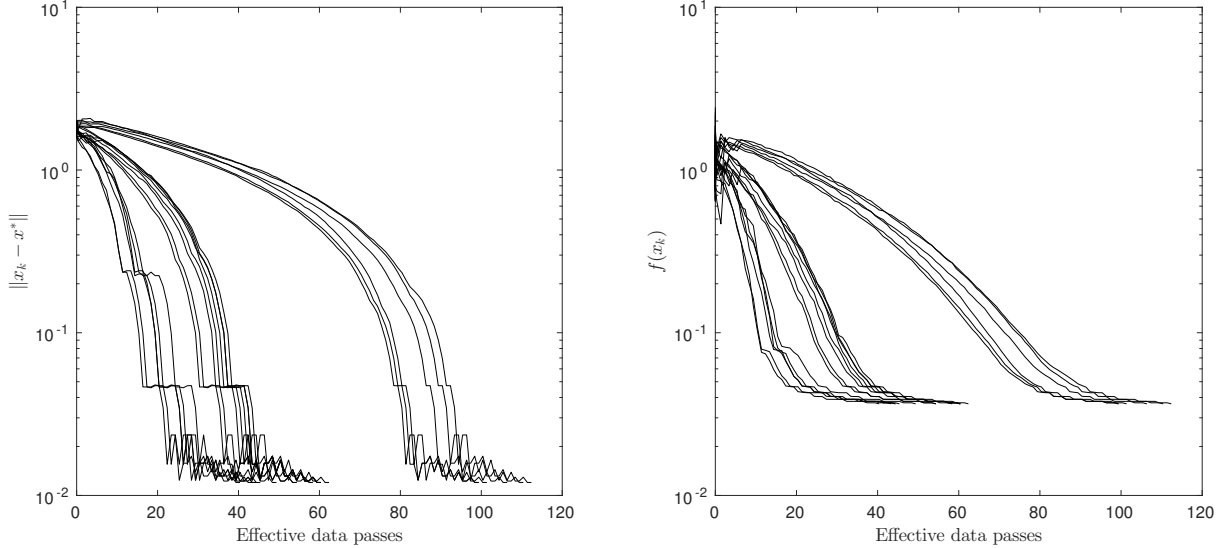
367 Given N , for each evaluation of f , we independently and randomly sample vectors (c_i, y_i) , $i = 1, \dots, N$
368 from the data set (c_i, y_i) , $i = 1, \dots, 10^7$ generated above. Note that we sample with replacement, following
369 the bootstrapping technique in statistics [19]. This allows the sample size to be larger than 10^7 . Then we
370 compute smoothing approximation $\tilde{f}(x, \mu, N)$ of the following function

$$371 \quad (4.4) \quad \tilde{f}(x, N) = \frac{1}{N} \sum_{i=1}^N [(\max(c_i^T x, 0) - y_i)^2] + \lambda \sum_{i=1}^n \log(1 + |(x)_i|)$$

372 by using smoothing functions for $\max(\cdot, 0)$ and $|\cdot|$.

373 For the initial point $x_0 = (0, \dots, 0)^T$ we performed 20 runs of Algorithm `mc_smoothing_search`. In
374 Figure 4.1, we show histories of the distance $\|x_k - x^*\|$ and the value $f(x_k)$.

FIG. 4.1. Histories of the distance $\|x_k - x^*\|$ and the value $f(x_k)$



375 Figure 4.1 illustrates several properties of the algorithm and the problem. At the end of the iteration,
 376 all of the iteration histories are very similar. The theory would lead one to expect similar histories if N
 377 is large. On the other hand, the initial value of N is large enough to cause considerable variation in f
 378 early in the iteration. This variation accounts for the differences in the histories. Finally, the iteration
 379 stagnates in the terminal phase when the differences from the iterates and x^* are roughly at the level of the
 380 regularization parameter. The reason for this is that the regularization term would dominate the error term
 381 when x is near x^* . While a smaller regularization would defer the stagnation, it would make it harder to
 382 capture the sparsity. Our choice of $\lambda = 10^{-2}$ captures the sparsity exactly. At each final iterate x_k , we have
 383 $(x_k)_i = (x^*)_i = 0$ for all $i \in S^c$, where $S = \{i \mid (x^*)_i \neq 0, i = 1, \dots, n\}$, the support set of x^* .

384 Figure 4.2 plots $\|x_k - x^*\|$ and $f(x_k)$ against the number of *Effective Data Passes* performed until the
 385 k -th iteration (EDP_k), which is defined as the number of samples made during the first k iterations divided
 386 by the data size 10^7 . Note that EDP_k increases with respect to k , but not strictly. The number of Effective
 387 Data Passes may stay unchanged during several successive iterations because the former changes only after
 388 a stencil failure, which happens only once in a few iterations. Consequently, a single number of Effective
 389 Data Passes can correspond to several values of $\|x_k - x^*\|$ or $f(x_k)$. That is why we observe vertical lines in
 390 Figure 4.2.

391 Figure 4.2 shows that the algorithm is capable of achieving considerable progress using very few Effective
 392 Data Passes. Both $\|x_k - x^*\|$ and $f(x_k)$ are reduced significantly even before one single Effective Data Pass
 393 is made. This shows the effectiveness of our sampling strategy, which increases the sample size steadily in
 394 course of the iterations. The stagnations in the final stage of the plots are seemingly more visible than in
 395 Figure 4.1. This is because the sample size increases rapidly, and hence the variations in $\|x_k - x^*\|$ or $f(x^*)$
 396 are less visible when plotted against EDP_k than plotted against k .

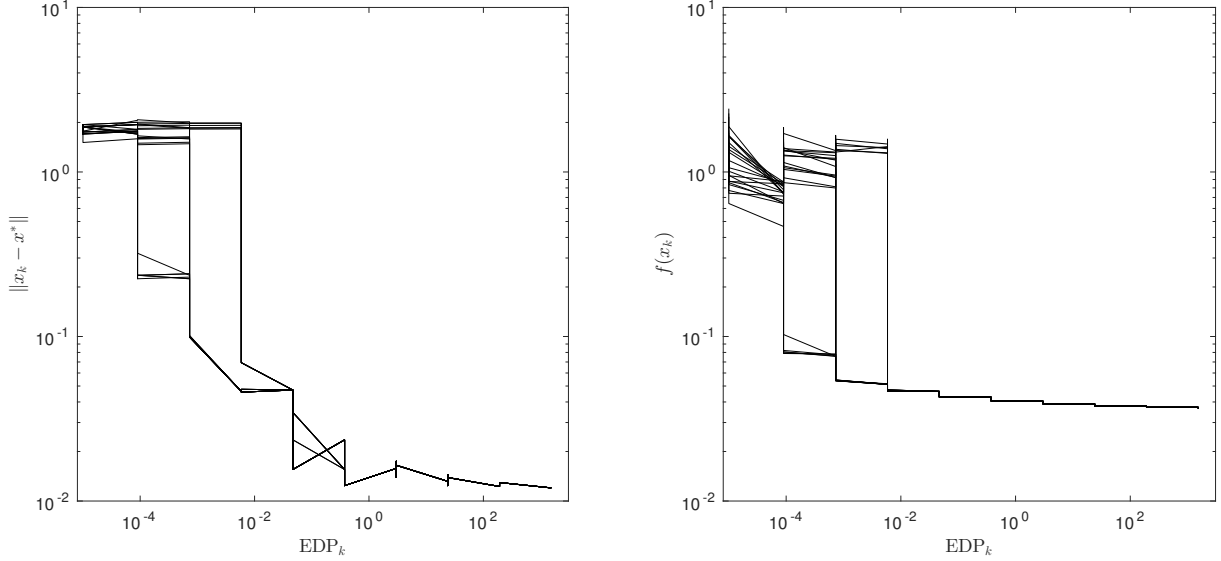
397 Note that we begin with $h = 0.5$ and terminate when $h \leq 10^{-3}$. By the structure of Algorithm
 398 `mc_smoothing_search`, the algorithmic parameters $\tau = 0.5$, $\gamma = 1.5$ and the initial values $N_0 = 100$,
 399 $\mu_0 = 0.1$, we know that after $t = 9$ iterations, $\tilde{h}_9 \leq 10^{-3}$, and hence $\tilde{\mu}_9 = \mu_0/2^{t\tau} = 0.1/2^{9 \times 0.5} = 0.044$ and
 400 $\tilde{N}_9 = 4^{t\gamma} N_0 = 4^{9 \times 1.5} 100 \approx 1.34\text{E}10$.

401 **4.2. Portfolio management.** Consider ν assets. Let $u \in R^\nu$ denote the random returns of them, and

402 (4.5)
$$r = E[u], \quad C = E[(u - r)(u - r)^T].$$

403 Here r is the vector of expected returns of the different assets, and C is the covariance matrix of the return
 404 on the assets in the portfolio. When r and C are known, as discussed in [44], the Markowitz mean-variance

FIG. 4.2. $\|x_k - x^*\|$ and $f(x_k)$ plotted against the number of Effective Data Passes EDP_k



405 model [36, 37] for portfolio selection can be formulated as

$$\begin{aligned}
 & \min_w \frac{1}{2} w^T C w - \eta r^T w \\
 & \text{s.t. } e^T w = 1 \\
 & \quad a \leq w \leq b,
 \end{aligned}
 \tag{4.6}$$

407 where w denotes the weights of the assets in the portfolio, $a \in R^\nu$ and $b \in R^\nu$ ($a \leq b$) are lower and upper
 408 bounds enforced on w , and η is a nonnegative parameter (called the risk aversion factor) to balance the
 409 conflicting aspects of minimizing the risk measured by $w^T C w$ and maximizing the expected return measured
 410 by $r^T w$. The Markowitz mean-variance model [36, 37] was first proposed and solved when the total return
 411 is known. The model captures the essence of two conflicting aspects in portfolio management; namely, the
 412 risk and the return.

413 The use of mean variance analysis in portfolio selection normally requires the knowledge of means,
 414 variances, and covariances of returns of all securities under consideration. However, in general, these data
 415 are not known exactly. Treating their estimates as if they were the exact parameters can lead to suboptimal
 416 portfolio choices.

417 The experiments reported in [22, 25, 30] show that, influenced by the sampling error, portfolios selected
 418 with the mean-variance model by Markowitz are not as efficient as an equally weighted portfolio. Other
 419 results [13, 39] show that the mean-variance model tends to magnify the errors associated with the estimates.
 420 In this section, we consider an optimal parameter selection model based on the Markowitz mean-variance
 421 model to find optimal parameters for portfolio selection.

For simplicity, here we only consider the case where C is positive definite and the feasible set $\{w \mid e^T w = 1, a \leq w \leq b\}$ is nonempty. Given a , b , and η , problem (4.6) has a unique solution w . In other words, w is uniquely defined by a , b , and η , the values of which will determine the quality of the portfolio selected by solving problem (4.6). A common measure for the quality is the Sharpe ratio [50]

$$\text{SR} = \frac{r^T w}{\sqrt{w^T C w}}.$$

422 The Sharpe ratio characterizes how well the return of an asset compensates the investor for the risk taken.
 423 In general, a strategy is better than others if its Sharpe ratio is higher.

424 In practice, a , b , and η are usually set by investors empirically according to their preferences. We
 425 consider selecting them by solving the two-stage optimization problem

$$\begin{aligned}
 & \max_{(a,b,\eta) \in \Omega} \frac{r^T w(a,b,\eta)}{\sqrt{w(a,b,\eta)^T C w(a,b,\eta)}} \\
 426 \quad (4.7) \quad & \text{where } w(a,b,\eta) = \underset{w}{\operatorname{argmin}} \frac{1}{2} w^T C w - \eta r^T w \\
 & \text{s.t. } e^T w = 1 \\
 & \quad a \leq w \leq b,
 \end{aligned}$$

where the feasible set

$$\Omega = [a, \bar{a}] \times [b, \bar{b}] \times [\eta, \bar{\eta}] \quad \text{with} \quad \underline{\eta} < \bar{\eta}$$

is given. The number of variables of the first level problem is

$$\#\{i | \underline{a}_i < \bar{a}_i, i = 1, \dots, \nu\} + \#\{i | \underline{b}_i < \bar{b}_i, i = 1, \dots, \nu\} + 1.$$

427 For example, if we choose

$$428 \quad (4.8) \quad \underline{a}_i = \bar{a}_i = 0, \text{ for } i \neq 1, \quad \underline{a}_1 = 0, \bar{a}_1 = 1 \quad \text{and} \quad \underline{b}_i = \bar{b}_i = 1, \text{ for } i \neq 2, \quad \underline{b}_2 = 0, \bar{b}_2 = 1,$$

429 then the number of variables of the first level problem is 3.

430 Finding optimal parameters a, b, η is a challenging problem. Since C is positive definite, the second level
 431 optimization problem has a unique solution. Hence $w(a, b, \eta)$ is well defined, and it is Lipschitz continuous
 432 with respect to (a, b, η) . However, $w(a, b, \eta)$ is not differentiable, and the covariance matrix C and the vector
 433 of expected return r cannot be computed directly in general. We will use the barrier method [43, Chapter
 434 19] to solve the second stage problem of (4.7) and define a smoothing function $w_\mu(a, b, \eta)$ [43]. In particular,
 435 we use Algorithm `mc_smoothing_search` to solve

$$\begin{aligned}
 & \max_{(a,b,\eta) \in \Omega} \frac{r^T w_\mu(a,b,\eta)}{\sqrt{w_\mu(a,b,\eta)^T C w_\mu(a,b,\eta)}} \\
 436 \quad (4.9) \quad & \text{where } w_\mu(a,b,\eta) = \underset{w}{\operatorname{argmin}} \frac{1}{2} w^T C w - \eta r^T w - \mu \sum_{i=1}^{\nu} \log(s_i) - \mu \sum_{i=1}^{\nu} \log(t_i) \\
 & \text{s.t. } e^T w = 1 \\
 & \quad w - a - s = 0 \\
 & \quad b - w - t = 0.
 \end{aligned}$$

437 In this section, we report numerical results that we obtained with Algorithm `mc_smoothing_search` for
 438 five standard data sets from the OR-Library [5], the SSE50 index, the CSI 100 index, and the CSI 300 index
 439 from Shanghai-Shenzhen stock market. The data are the weekly or daily prices of the component stocks for
 440 the eight stock market indices drawn from different countries. See Table 4.1 for the description of the data
 441 sets. The ν and T columns are the number of the component assets included in the index and the number
 442 of the observations for the assets, respectively.

443 We report on two experiments: randomly generated problems in §4.2.1, which use the mean and the
 444 covariance matrix generated from the real data in Table 4.1 and rolling window procedures for out-of-sample
 445 comparison in §4.2.2, which use the stock prices to generate the returns of assets and the covariance matrix
 446 by Monte Carlo simulation.

4.2.1. Randomly generated problems. We choose the following parameters as input data of Algo-
 rithm `mc_smoothing_search`:

$$h = 0.5, \quad \mu = 0.1, \quad N = 100, \quad \tau = 0.5, \quad \gamma = 1.5.$$

447 We choose the feasible set X as in (4.8).

TABLE 4.1
Description of the eight real data sets

Data set	ν	Location	Index	T	Description
data 1	31	Hong Kong	Hang Seng	291	weekly prices from 1992 to 1997
data 2	85	Germany	DAX 100	291	weekly prices from 1992 to 1997
data 3	89	UK	FTSE 100	291	weekly prices from 1992 to 1997
data 4	98	USA	S&P 100	291	weekly prices from 1992 to 1997
data 5	225	Japan	Nikkei 225	291	weekly prices from 1992 to 1997
SSE50	50	China	SSE50	501	daily prices from 2013 to 2015
CSI100	100	China	CSI 100	501	daily prices from 2013 to 2015
CSI300	300	China	CSI 300	401	daily prices from 2011 to 2013

448 For all tests, we terminated the algorithm if the stencil size is less than 10^{-2} .

For each data set in Table 4.1, we first calculate the average $\hat{r} \in R^\nu$ and the covariance matrix $\hat{C} \in R^{\nu \times \nu}$ for the returns of the assets. Given a sample size N , we generate i.i.d. random vectors $u_i \in R^\nu$, $i = 1, \dots, N$ normally distributed with mean \hat{r} and covariance matrix \hat{C} , that is

$$u_i = \hat{r} + \hat{C}^{\frac{1}{2}} \text{randn}(\nu, 1), \quad i = 1, 2, \dots, N,$$

and then take r^N to be the sample average of u_i , $i = 1, 2, \dots, N$, and C^N to be the sample covariance matrix

$$r^N = \frac{1}{N} \sum_{i=1}^N u_i \quad \text{and} \quad C^N = \frac{1}{N} \sum_{i=1}^N (u_i - r^N)(u_i - r^N)^T.$$

449 Then we compute the smoothing approximation for the problem of minimizing the negative Sharpe ratio
450 which is given by (4.9) with $r = r^N$, $C = C^N$, and $x = (a_1, b_2, \eta)$.

451 For each data set in Table 4.1, we use Algorithm `mc_smoothing_search` to solve problem (4.7) with the
452 starting point $(a_1, b_2, \eta) = (0, 1, 0.5)$. Table 4.2 presents the results. From Table 4.2, we can see the optimal
453 value of objective function (Opt. sharpe ratio) at the final iteration is bigger than the value of objective
454 function at the point $w = e/\nu$, which is a feasible point of problem (4.6) with $a = 0$ and $b = e$.

TABLE 4.2
Numerical results for the portfolio management problem with randomly generated data

data set	data1	data2	data3	data4	data5	SSE50	CSI100	CSI300
lower bound a_1	4.69E-1	4.22E-1	6.56E-1	1.25E-1	1.00E00	2.00E-2	1.00E-2	7.50E-1
upper bound b_2	8.13E-1	8.75E-1	9.69E-1	9.69E-1	7.12E-1	1.14E-1	3.85E-1	7.50E-1
risk aversion η	9.69E-1	5.31E-1	3.43E-1	4.38E-1	8.59E-1	1.25E-1	6.25E-1	6.41E-1
opt. sharpe ratio	1.57E-1	2.85E-1	2.51E-1	2.47E-1	9.76E-2	3.92E-1	2.74E-1	7.98E-2
sharpe ratio e/ν	1.04E-1	9.15E-2	1.53E-1	1.99E-1	-4.90E-2	2.36E-1	2.65E-1	-9.05E-2

4.2.2. Problems with rolling window procedures. For a given data set, assuming that the obser-
vations of the stock prices are $\{P_{i,t} : 1 \leq i \leq \nu, 1 \leq t \leq T\}$, we can compute the (logarithmic) returns of
the stocks:

$$r_{i,t} = \log \frac{P_{i,t+1}}{P_{i,t}}, \quad i = 1, \dots, \nu, \quad t = 1, \dots, T - 1.$$

455 For the purpose of numerical comparisons, we partition the data set into two subsets: a training set and a
456 testing set. The training set, called in-sample set, consists of the first half of the data set and is used to
457 compute an optimal parameter x^* and the corresponding optimal portfolio selection $w(x^*)$. The testing set,

458 called out-of-sample, consists of the second half of the data set and is used to test how well the optimal
 459 parameter x^* and the corresponding optimal portfolio selection $w(x^*)$.

More exactly, for stock i with $i = 1, \dots, \nu$, we can use the training set to compute the in-sample expectation and the standard deviation by

$$\bar{\mu}_i = \frac{1}{M} \sum_{t=1}^M r_{i,t} \quad \text{and} \quad \bar{\sigma}_i = \sqrt{\frac{1}{M} \sum_{t=1}^M (r_{i,t} - \bar{\mu}_i)^2},$$

respectively, where $M = (T - 1)/2$. As is standard in finance [24], we simulate the out-of-sample prices as follows. Let N be the sample size. Then at the j -th simulation ($1 \leq j \leq N$), for $M + 1 \leq t \leq T - 1$, if the price $S_{i-1,t}^{(j)}$ of stock i at an out-of-sample time $t - 1$ is known, the price $S_{i,t}^{(j)}$ of this stock at time t is generated by

$$S_{i,t}^{(j)} = S_{i,t-1}^{(j)} \exp(\bar{\mu}_i + \bar{\sigma}_i Z),$$

where $S_{i,M}^{(j)} = P_{i,M}$ for all $1 \leq j \leq N$ and Z is randomly produced by the standard normal distribution $N(0, 1)$. In a similar way, we can calculate the (logarithmic) returns by this simulation

$$r_{i,t}^{(j)} = \log \frac{S_{i,t+1}^{(j)}}{S_{i,t}^{(j)}}, \quad t = M + 1, \dots, T - 1.$$

For $t = M + 1, \dots, T - 1$, denote the column vector $r_t^{(j)}$ with its i -th component being $r_{i,t}^{(j)}$ and its average vector $\bar{r}_t = \frac{1}{N} \sum_{j=1}^N r_t^{(j)}$, the sample mean r^N and the sample variance C^N of the out-of-sample can be computed by

$$r^N = \frac{1}{M} \sum_{t=M+1}^{T-1} \bar{r}_t \quad \text{and} \quad C^N = \frac{1}{M} \sum_{t=M+1}^{T-1} (\bar{r}_t - r^N)(\bar{r}_t - r^N)^T.$$

460 Then we solve problem (4.9) with the sample mean r^N and the sample variance C^N to obtain the optimal
 461 parameter x^* and the corresponding optimal portfolio selection $w(x^*)$ by Algorithm `mc_smoothing_search`.

We choose the following parameters as input data of Algorithm `mc_smoothing_search`:

$$h = 0.5, \quad \mu = 0.1, \quad N = 10, \quad \tau = 0.5, \quad \gamma = 1.5.$$

462 We choose the feasible set X as in (4.8). For all tests, we choose the starting point $(a_1, b_2, \eta) = (0, 1, 0.5)$
 463 and terminated the algorithm when the sample size N gets larger than 10^5 .

To evaluate the quality of the optimal portfolio selection $w(x^*)$, we shall make use of the real out-of-sample data. We denote by r^{out} and C^{out} the mean and variance of the real returns of the out-of-sample set; namely,

$$r^{out} = \frac{1}{M} \sum_{t=M+1}^{T-1} r_t \quad \text{and} \quad C^{out} = \frac{1}{M} \sum_{t=M+1}^{T-1} (r_t - r^{out})(r_t - r^{out})^T,$$

where r_t is the vector formed by the stock prices $r_{i,t}$ ($i = 1, \dots, n$) for $t = M + 1, \dots, T - 1$. Then we can calculate the Sharpe ratio of the optimal solution $w(x^*)$ by using r^{out} and C^{out} as

$$SR^* = \frac{(r^{out})^T w(x^*)}{\sqrt{w(x^*)^T C^{out} w(x^*)}}.$$

464 In Table 4.3, for all the eight data sets, we list the optimal values of a_1 , b_2 , η achieved by Algorithm
 465 `mc_smoothing_search`, and the corresponding SR^* . For comparison, we also list the Sharpe ratio of the
 466 average strategy (namely, taking $1/\nu$ portion of each portfolio) using r^{out} and C^{out} . From Table 4.3, we can
 467 see that using Algorithm `mc_smoothing_search` to solve problem (4.7) can provide a portfolio strategy with
 468 higher Sharpe ratio than the average strategy for all data sets.

TABLE 4.3
Numerical results for the portfolio management problem with rolling window procedures

data set	data1	data2	data3	data4	data5	SSE50	CSI100	CSI300
lower bound a_1	1.00E-3	1.18E-2	1.12E-2	1.02E-1	4.40E-2	2.00E-2	1.00E-2	3.33E-3
upper bound b_2	3.14E-1	2.31E-1	6.36E-1	4.15E-2	1.29E-2	5.83E-1	2.60E-1	2.53E-1
risk aversion η	2.81E-1	9.38E-1	6.25E-1	6.25E-2	1.00E00	7.19E-1	2.50E-1	7.50E-1
sharpe ratio SR^*	3.35E-1	2.36E-1	3.72E-1	5.12E-1	2.19E-1	2.71E-1	4.33E-1	2.19E-1
sharpe ratio e/ν	1.57E-1	2.10E-1	2.79E-1	3.44E-1	-3.85E-2	2.36E-1	2.65E-1	3.18E-3

469 **5. Conclusions.** In this paper we propose a smoothing direct search algorithm with Monte Carlo sim-
470 ulation **Algorithm mc_smoothing_search** for the constrained nonsmooth nonconvex optimization problem
471 (1.2), where the objective function value $f(x)$ cannot be computed directly, but are approximated by Monte
472 Carlo simulation. This algorithm updates the stencil size h , smoothing parameter μ and the sample size N
473 simultaneously with the rate $h/\mu \rightarrow 0$, and $(h\sqrt{N})^{-1} \rightarrow 0$. We prove that any accumulation point of the
474 sequence generated by the algorithm satisfies the first order optimality condition $\tilde{\chi}(x) = 0$ with probability
475 one, where $\tilde{\chi}(x)$ is defined by (2.22). We report on a set of numerical experiments which illustrate the analysis
476 and show that **Algorithm mc_smoothing_search** is an effective method for minimizing nonsmooth functions
477 whose function values cannot be computed directly but are approximated by Monte Carlo simulation.

478 **Appendix A. Smoothing functions.**

479 We give an example of smoothing functions to explain Assumption 2.1. Let $f(x) = 2 \max(0, p(x))$, where
480 $p : R^n \rightarrow R$ is twice continuously differentiable with

481
$$\|\nabla p(x)\nabla p(x)^T\| \leq \Gamma.$$

482 We use the smoothing function

483 (A.1)
$$\hat{f}(x, \mu) = p(x) + \sqrt{p(x)^2 + 4\mu^2},$$

484 and $V = \{e_1, \dots, e_n\}$, the unit coordinate directions in R^n .

Clearly part (i) of Assumption 2.1 holds with $c_1 = 2$ and $c_2 = 0$, since

$$|f(x) - \hat{f}(x, \mu)| \leq 2\mu.$$

485 Now we consider part (ii) of Assumption 2.1. The Clarke subgradient has the form

486 (A.2)
$$\partial f(x) = 2 \begin{cases} \nabla p(x) & \text{if } p(x) > 0 \\ \mathbf{0} & \text{if } p(x) < 0 \\ [0, 1]\nabla p(x) & \text{if } p(x) = 0 \end{cases}$$

487 and the gradient of the smoothing function is

488
$$\nabla \hat{f}(x, \mu) = \left(1 + \frac{p(x)}{\sqrt{p(x)^2 + 4\mu^2}}\right) \nabla p(x).$$

Hence, we have

$$\|\nabla \hat{f}(x, \mu)\| \leq 2\|\nabla p(x)\|.$$

489 If p is Lipschitz continuously differentiable on a convex and compact set Ω , then there is an Υ such that
490 $\|\nabla \hat{f}(x, \mu)\| \leq \Upsilon$ on Ω .

It is easy to see that for $p(x) \neq 0$, f is differentiable at x and

$$\partial f(x) = \nabla f(x) = \text{con}\{v \mid \nabla \hat{f}(x_k, \mu_k) \rightarrow v, \text{ for } x_k \rightarrow x, \mu_k \downarrow 0\}.$$

For $p(x) = 0$, since $0 \leq 1 + \frac{p(x)}{\sqrt{p(x)^2 + 4\mu^2}} \leq 2$, we have

$$\text{con}\{v \mid \nabla \hat{f}(x_k, \mu_k) \rightarrow v, \text{ for } x_k \rightarrow x, \mu_k \downarrow 0\} \subseteq \partial f(x).$$

Now, let $\mu_k^2 = (1 - h^2)p(x_k)^2/(4h^2)$ for some $h \in (0, 1]$. Then for $x_k \rightarrow x$ with $p(x_k) \downarrow 0$, we have $\mu_k \downarrow 0$ and

$$\nabla \hat{f}(x_k, \mu_k) = (1 + h)\nabla p(x_k) \rightarrow (1 + h)\nabla p(x),$$

and for $x_k \rightarrow x$ with $p(x_k) \uparrow 0$, we have $\mu_k \downarrow 0$ and

$$\nabla \hat{f}(x_k, \mu_k) = (1 - h)\nabla p(x_k) \rightarrow (1 - h)\nabla p(x).$$

Moreover, if we take $\mu_k = \sqrt{|p(x_k)|}$, then for $x_k \rightarrow x$, we have $p(x_k) \rightarrow 0$, $\mu_k \downarrow 0$ and

$$\nabla \hat{f}(x_k, \mu_k) = \left(1 + \frac{p(x_k)}{\sqrt{p(x_k)^2 + 4|p(x_k)|}}\right) \nabla p(x_k) \rightarrow \nabla p(x).$$

Hence we find that for $p(x) = 0$,

$$\partial f(x) = [0, 2]\nabla p(x) = \text{con}\{v \mid \nabla \hat{f}(x_k, \mu_k) \rightarrow v, \text{ for } x_k \rightarrow x, \mu_k \downarrow 0\}.$$

491 Finally, we consider part (iii) of Assumption 2.1. Since

$$492 \quad \nabla^2 \hat{f}(x, \mu) = \left(1 + \frac{p(x)}{\sqrt{p(x)^2 + 4\mu^2}}\right) \nabla^2 p(x) + \frac{4\mu^2}{(p(x)^2 + 4\mu^2)^{\frac{3}{2}}} \nabla p(x) \nabla p(x)^T,$$

493 we have

$$494 \quad \|\nabla^2 \hat{f}(x, \mu)\| \leq \frac{1}{2\mu} \|\nabla p(x) \nabla p(x)^T\| + 2\|\nabla^2 p(x)\| \leq \frac{\Gamma}{2\mu} + 2\|\nabla^2 p(x)\|,$$

495 which implies that part (iii) of Assumption 2.1 holds.

496

497 A smoothing function of $|p(x)|$ can be defined by using the relation $|p(x)| = \max(0, p(x)) + \max(0, -p(x))$
 498 and a smoothing function of $\max(0, p(x))$. For example, using (A.1), we can have a smoothing function
 499 $\sqrt{(p(x)^2 + 4\mu^2)}$ for $|p(x)|$.

500 There is a detailed discussion of smoothing functions in [10].

501 Appendix B. The proof of (3.4).

502 In this appendix, we will show the existence of $\delta \in (0, 1)$, $c_F > 0$, and $\bar{N} > 0$ that fulfill (3.4). To this
 503 end, we have to study the uniform convergence rate of the empirical mean (3.3), which has been investigated
 504 in [26, 59] under certain conditions. In particular, if $\phi(\xi, x) = \mathbf{1}(\xi \leq x)$, where $\xi, x \in \mathbb{R}$, and $\mathbf{1}(E)$ is
 505 the indicator function of an event E , then inequality (3.4) is indeed satisfiable with $p = 1/2$ (see [18, 38]).
 506 However, in general, one can only achieve (3.4) for $p < 1/2$. Our argument here is essentially an extension
 507 of the discussions in [59, Section 3].

508 Before the proof, we recall that $\phi(\xi, x)$ is sub-exponential for each $x \in X$. As defined in [56, Definition
 509 5.13], a real-value random variable ζ is called a sub-exponential¹ random variable if

$$510 \quad (\text{B.1}) \quad \sup_{p \geq 1} p^{-1} E(|\zeta|^p)^{1/p} < \infty.$$

511 The quantity on the left-hand side of (B.1), often denoted by $\|\zeta\|_{\psi_1}$, is called the sub-exponential norm
 512 of ζ . According to [17, Theorem 3.14], ζ is sub-exponential as per (B.1) if and only if its *moment generating*

¹ Note that there is another widely used but completely different concept of sub-exponentiality in probability theory, which refers to a certain heavy-tail behavior of distributions as detailed in [21, 54]. The sub-exponentiality defined by (B.1) is commonly found in areas like machine learning and data analysis (see [17, Section 3.1.2] and [57, Chapter 2]), and is a slight generalization of the pre-Gaussianity defined in [7, Chapter 1].

513 function $G(t) \equiv E(e^{\zeta t})$ is finite in a neighbourhood of zero (see also [56, inequality (5.16)] and [57, Theorem
514 2.2]). As elaborated in [56, inequality (5.14)] and [57, Theorem 2.2], assuming sub-exponentiality is also
515 equivalent to requiring that the tail probability $P(|\zeta| > t)$ decays exponentially or faster, which is certainly
516 not a trivial condition. However, it encompasses a large class of distributions that are interesting in practice,
517 including all the distributions with bounded support sets, the normal distribution, Gamma distribution,
518 Weibull distribution, Poisson distribution, geometric distribution, and any Lipschitz continuous functions of
519 random variables following such distributions (see standard statistics textbooks like [40] for the definitions
520 and moment generating functions of the named distributions).

521 An important property of sub-exponential random variables is a Bernstein-type inequality presented
522 in [56, Proposition 5.16]: if $\zeta_1, \zeta_2, \dots, \zeta_N$ are independent sub-exponential random variables with zero
523 mean, and $\sigma = \max_{1 \leq i \leq N} \|\zeta_i\|_{\psi_1}$, then

$$524 \quad (\text{B.2}) \quad \text{Prob} \left(\frac{1}{N} \left| \sum_{i=1}^N \zeta_i \right| \geq t \right) \leq 2 \exp \left(-cN \min \left\{ \frac{t^2}{\sigma^2}, \frac{t}{\sigma} \right\} \right) \quad \text{for each } t \geq 0,$$

525 where c is an *absolute constant*. Therefore,

$$526 \quad (\text{B.3}) \quad \text{Prob} \left(\frac{1}{N} \left| \sum_{i=1}^N \zeta_i \right| \geq t \right) \leq 2 \exp \left(-\frac{cNt^2}{\sigma^2} \right) \quad \text{when } 0 \leq t \leq \sigma.$$

527 More general forms of inequalities (B.2) and (B.3) can be found in [9, Theorem 1.2.7] and [17, Corollary
528 3.17]. Note that (B.2) still holds if we replace σ with any number larger than $\max_{1 \leq i \leq N} \|\zeta_i\|_{\psi_1}$, for the
529 right-hand side of (B.2) is decreasing with respect to σ . Consequently, inequality (B.3) is valid as long as

$$530 \quad \sigma \geq \max_{1 \leq i \leq N} \|\zeta_i\|_{\psi_1}.$$

531 In Section 3, by saying that the m -dimensional random vector $\phi(\xi, x)$ is sub-exponential, we mean that
532 each entry of $\phi(\xi, x)$ is a sub-exponential random variable (not necessarily independent of each other). Then,
533 according to [56, Remark 5.18], each entry of $F(x) - \phi(\xi, x)$ is a sub-exponential random variable with zero
534 mean, since $F(x) = E[\phi(\xi, x)]$ as stated in (3.2).

535 Now we give the proof of (3.4).

536 *Proof.* Consider an arbitrary point $y \in X$. Let $\left(F(y) - \phi(\xi, y) \right)_i$ be the i -th entry of the random vector
537 $F(y) - \phi(\xi, y)$, and

$$538 \quad \sigma(y) = 1 + \max_{1 \leq i \leq m} \left\| \left(F(y) - \phi(\xi, y) \right)_i \right\|_{\psi_1}.$$

539 For each $i \in \{1, 2, \dots, m\}$, since $\sigma(y) \geq \left\| \left(F(y) - \phi(\xi, y) \right)_i \right\|_{\psi_1}$, we can invoke the Bernstein-type bound (B.3)
540 to obtain

$$541 \quad (\text{B.4}) \quad \text{Prob} \left(\left| \left(F(y) - \tilde{F}^N(y) \right)_i \right| \geq t \right) = \text{Prob} \left(\frac{1}{N} \left| \sum_{\ell=1}^N \left(F(y) - \phi(\xi_\ell, y) \right)_i \right| \geq t \right) \leq 2 \exp \left[-\frac{cNt^2}{\sigma^2(y)} \right]$$

542 whenever $t \in [0, 1] \subset [0, \sigma(y)]$. Hence

$$543 \quad (\text{B.5}) \quad \text{Prob} \left(\|F(y) - \tilde{F}^N(y)\| \geq \epsilon \right) \leq \text{Prob} \left(\max_{1 \leq i \leq m} \left| \left(F(y) - \tilde{F}^N(y) \right)_i \right| \geq \frac{\epsilon}{\sqrt{m}} \right) \leq 2m \exp \left[-\frac{cN\epsilon^2}{m\sigma^2(y)} \right]$$

544 for each $\epsilon \in [0, 1]$. This gives us the point-wise convergence rate of \tilde{F}^N . In the following, we will extend this
545 to obtain an estimation for the uniform convergence rate. The key is to exploit the Lipschitz continuity of
546 ϕ and the boundedness of X .

547 Since $\phi(\xi, x)$ is L -Lipschitz with respect to x for a constant L independent of ξ , both F and \tilde{F}^N are
548 L -Lipschitz continuous. Let D be the diameter of X , which is finite since X is bounded. Then there exists a
549 set $\{y_j\}_{j=1}^K \subseteq X$ such that

$$550 \quad (\text{B.6}) \quad K \leq \left\lceil \frac{\sqrt{n}D}{\epsilon/(4L)} \right\rceil^n \quad \text{and} \quad X \subseteq \bigcup_{j=1}^K B\left(y_j, \frac{\epsilon}{4L}\right).$$

551 For each $x \in X$, let j_x be an integer in $\{1, 2, \dots, K\}$ such that $x \in B(y_{j_x}, \epsilon/(4L))$. Then

$$\begin{aligned} 552 \quad \|F(x) - \tilde{F}^N(x)\| &\leq \|F(x) - F(y_{j_x})\| + \|\tilde{F}^N(x) - \tilde{F}^N(y_{j_x})\| + \|F(y_{j_x}) - \tilde{F}^N(y_{j_x})\| \\ &\leq L\|x - y_{j_x}\| + L\|x - y_{j_x}\| + \max_{1 \leq j \leq K} \|F(y_j) - \tilde{F}^N(y_j)\| \\ &\leq \frac{\epsilon}{2} + \max_{1 \leq j \leq K} \|F(y_j) - \tilde{F}^N(y_j)\|. \end{aligned}$$

553 This, together with (B.5) and (B.6), tells us that

$$\begin{aligned} 554 \quad \text{Prob}\left(\sup_{x \in X} \|F(x) - \tilde{F}^N(x)\| \geq \epsilon\right) &\leq \text{Prob}\left(\max_{1 \leq j \leq K} \|F(y_j) - \tilde{F}^N(y_j)\| \geq \frac{\epsilon}{2}\right) \\ &\leq \sum_{j=1}^K \left\{ 2m \exp\left[-\frac{cN\epsilon^2}{4m\sigma^2(y_j)}\right] \right\} \\ &\leq 2m \left\lceil \frac{\sqrt{n}D}{\epsilon/(2L)} \right\rceil^n \exp(-c' N\epsilon^2) \end{aligned}$$

for each $\epsilon \in [0, 1]$, where

$$c' \equiv \min_{1 \leq j \leq K} \frac{c}{4m\sigma^2(y_j)} > 0.$$

555 Setting $\epsilon = 1/N^p$, we obtain

$$556 \quad (\text{B.7}) \quad \text{Prob}\left(\sup_{x \in X} \|F(x) - \tilde{F}^N(x)\| \geq \frac{1}{N^p}\right) \leq 2m [2\sqrt{n}LDN^p]^n \exp(-c' N^{1-2p}).$$

557 Given $p \in (0, 1/2)$, we can choose \bar{N} large enough (depending on m, n, L, D, c' , and p) such that the
558 right-hand side of (B.7) is at most $1/2$ for each $N \geq \bar{N}$. Consequently,

$$559 \quad (\text{B.8}) \quad \text{Prob}\left(\sup_{x \in X} \|F(x) - \tilde{F}^N(x)\| \geq \frac{1}{N^p}\right) \leq \frac{1}{2}$$

560 for each $N \geq \bar{N}$. In other words, this \bar{N} fulfills (3.4) together with $\delta = 1/2$ and $c_F = 1$. \square

561 **Acknowledgement.** The authors are grateful to three referees and the associate editor for their very
562 helpful comments and suggestions.

563

REFERENCES

- 564 [1] J. G. ALTONJI, H. ICHIMURA, AND T. OTSU, *Estimating derivatives in nonseparable models with limited dependent vari-*
565 *ables*, *Econometrica*, 80 (2012), pp. 1701–1719.
566 [2] C. AUDET AND J. E. DENNIS, *Analysis of generalized pattern searches*, *SIAM J. Optim.*, 13 (2003), pp. 889–903.
567 [3] C. AUDET AND J. E. DENNIS, *Mesh adaptive direct search algorithms for constrained optimization*, *SIAM J. Optim.*, 17
568 (2006), pp. 188–217.
569 [4] C. AUDET AND J. E. DENNIS, *A progressive barrier for derivative-free nonlinear programming*, *SIAM J. Optim.*, 20 (2009),
570 pp. 445–472.
571 [5] J. E. BEASLEY, *OR-Library: distributing test problems by electronic mail*, *J. Oper. Res. Soc.*, (1990), pp. 1069–1072.
572 [6] R. W. BLUNDELL AND J. L. POWELL, *Censored regression quantiles with endogenous regressors*, *J. Econometrics*, 141
573 (2007), pp. 65–83.

- 574 [7] V. V. BULDYGIN AND I. U. V. KOZACHENKO, *Metric characterization of random variables and random processes*, American
575 Mathematical Society, Providence, 2000.
- 576 [8] J. V. BURKE, X. CHEN, AND H. SUN, *Compute Clarke subgradients of expectation of measurable composite affine max*
577 *function via smoothing*, 2017. manuscript.
- 578 [9] D. CHAFAÏ, O. GUÉDON, G. LECUÉ, AND A. PAJOR, *Interactions between compressed sensing random matrices and high*
579 *dimensional geometry*, Société Mathématique de France, 2012, <http://lecueguillaume.github.io/assets/CSbook.pdf>.
- 580 [10] X. CHEN, *Smoothing methods for nonsmooth, nonconvex minimization*, Math. Program., 134 (2012), pp. 71–99.
- 581 [11] X. CHEN AND C. T. KELLEY, *Optimization with hidden constraints and embedded Monte Carlo computations*, Optim.
582 Eng., 17 (2016), pp. 157–175.
- 583 [12] X. CHEN, T. K. PONG, AND R. J.-B. WETS, *Two-stage stochastic variational inequalities: an ERM-solution procedure*,
584 Math. Program., 165 (2017), pp. 71–112.
- 585 [13] V. K. CHOPRA, *Mean-variance revisited: near-optimal portfolios and sensitivity to input variations*, J. Investing, 2 (1993),
586 pp. 51–59.
- 587 [14] F. H. CLARKE, *Optimization and Nonsmooth Analysis*, no. 5 in Classics in Applied Mathematics, SIAM, Philadelphia,
588 1990.
- 589 [15] A. R. CONN, K. SCHEINBERG, AND L. N. VICENTE, *Introduction to Derivative-Free Optimization*, MPS-SIAM Series on
590 Optimization, SIAM, Philadelphia, 2009.
- 591 [16] J. E. DENNIS AND V. TORCZON, *Direct search methods on parallel machines*, SIAM J. Optim., 1 (1991), pp. 448 – 474.
- 592 [17] J. DUCHI, *Lecture notes for Statistics 311/Electrical Engineering 377: Information theory and statistics*, [https://stanford.](https://stanford.edu/class/stats311/Lectures/full_notes.pdf)
593 [edu/class/stats311/Lectures/full_notes.pdf](https://stanford.edu/class/stats311/Lectures/full_notes.pdf).
- 594 [18] A. DVORETZKY, J. KIEFER, AND J. WOLFOWITZ, *Asymptotic minimax character of the sample distribution function and*
595 *of the classical multinomial estimator*, Ann. Math. Statist., 27 (1956), pp. 642–669.
- 596 [19] B. EFRON AND R. J. TIBSHIRANI, *An Introduction to the Bootstrap*, CRC Press, New York, 1994.
- 597 [20] F. FACCHINEI AND J.-S. PANG, *Finite-dimensional Variational Inequalities and Complementarity Problems*, Springer-
598 Verlag, New York, 2003.
- 599 [21] S. FOSS, D. KORSHUNOV, AND S. ZACHARY, *An Introduction to Heavy-Tailed and Subexponential Distributions*, Springer,
600 New York, 2011.
- 601 [22] G. M. FRANKFURTER, H. E. PHILLIPS, AND J. P. SEAGLE, *Portfolio selection: the effects of uncertain means, variances,*
602 *and covariances*, J. Financ. Quantit. Anal., 6 (1971), pp. 1251–1262.
- 603 [23] R. GARMANJANI AND L. N. VICENTE, *Smoothing and worst-case complexity for direct-search methods in nonsmooth*
604 *optimization*, IMA J. Numer. Anal., 33 (2012), pp. 1008–1028.
- 605 [24] J. C. HULL, *Options, Futures, and Other Derivatives*, Pearson, New York, 1988.
- 606 [25] P. JORION, *Bayes-Stein estimation for portfolio analysis*, J. Financ. Quantit. Anal., 21 (1986), pp. 279–292.
- 607 [26] R. L. KARANDIKAR AND M. VIDYASAGAR, *Rates of uniform convergence of empirical means with mixing processes*, Statist.
608 Probab. Lett., 58 (2002), pp. 297–307.
- 609 [27] C. T. KELLEY, *Iterative Methods for Optimization*, no. 18 in Frontiers in Applied Mathematics, SIAM, Philadelphia, 1999.
- 610 [28] C. T. KELLEY, *Implicit Filtering*, no. 23 in Software Environments and Tools, SIAM, Philadelphia, 2011.
- 611 [29] S. KIM, R. PASUPATHY, AND S. G. HENDERSON, *A guide to sample average approximation*, in Handbook of Simulation
612 Optimization, M. C. Fu, ed., Springer New York, New York, 2015, pp. 207–243.
- 613 [30] R. W. KLEIN AND V. S. BAWA, *The effect of estimation risk on optimal portfolio choice*, J. Financ. Econom., 3 (1976),
614 pp. 215–231.
- 615 [31] T. G. KOLDA, R. M. LEWIS, AND V. TORCZON, *Stationarity results for generating set search for linearly constrained*
616 *optimization*, SIAM J. Optim., 17 (2006), pp. 943–968.
- 617 [32] T. G. KOLDA, R. M. LEWIS, AND V. J. TORCZON, *Optimization by direct search: New perspectives on some classical and*
618 *modern methods*, SIAM Review, 45 (2003), pp. 385–482.
- 619 [33] H. J. KUSHNER AND G. G. YIN, *Stochastic Approximation and Recursive Algorithms and Applications*, Springer-Verlag,
620 New York, 2003.
- 621 [34] G. LAN, *An optimal method for stochastic composite optimization*, Math. Program., 133 (2012), pp. 365–397.
- 622 [35] Y.-F. LIU, S. MA, Y.-H. DAI, AND S. ZHANG, *A smoothing SQP framework for a class of composite l_q minimization over*
623 *polyhedron*, Math. Program., 158 (2016), pp. 467–500.
- 624 [36] H. MARKOWITZ, *Portfolio selection*, J. Finance, 7 (1952), pp. 77–91.
- 625 [37] H. MARKOWITZ, *The optimization of a quadratic function subject to linear constraints*, Naval Research Logistics Quarterly,
626 3 (1956), pp. 111–133.
- 627 [38] P. MASSART, *The tight constant in the Dvoretzky-Kiefer-Wolfowitz inequality*, Ann. Probab., 18 (1990), pp. 1269–1283.
- 628 [39] R. O. MICHAUD, *The Markowitz optimization enigma: Is optimized optimal?*, Financial Analysts J., 45 (1989), pp. 43–54.
- 629 [40] A. M. MOOD, *Introduction to the Theory of Statistics.*, McGraw-Hill, New York, 1950.
- 630 [41] Y. NESTEROV, *Introductory Lectures on Convex Optimization: A Basic Course*, vol. 87 of Applied Optimization, Kluwer
631 Academic Publishers, London, 2004.
- 632 [42] Y. NESTEROV AND V. SPOKOINY, *Random gradient-free minimization of convex functions*, Found. Comput. Math., 17
633 (2017), pp. 527–566.
- 634 [43] J. NOCEDAL AND S. WRIGHT, *Numerical Optimization*, Springer Series in Operations Research, Springer, New York,
635 second ed., 2006.
- 636 [44] J.-S. PANG, *A new and efficient algorithm for a class of portfolio selection problems*, Oper. Res., 28 (1980), pp. 754–767.
- 637 [45] R. PASUPATHY, *On choosing parameters in retrospective-approximation algorithms for stochastic root finding and simu-*
638 *lation optimization*, Oper. Res., 58 (2010), pp. 889–901.
- 639 [46] R. PASUPATHY, P. W. GLYNN, G. S. G., AND F. S. HAHEMI, *On sampling rates in simulation-based recursions*, Optimiza-

- 640 tion Online, (2017).
- 641 [47] M. PORCELLI AND P. L. TOINT, *BFO, a trainable derivative-free Brute Force Optimizer for nonlinear bound-constrained*
642 *optimization and equilibrium computations with continuous and discrete variables*, ACM Trans. Math. Software, (to
643 appear).
- 644 [48] S. SANKARAN, C. AUDET, AND A. L. MARSDEN, *A method for stochastic constrained optimization using derivative-free*
645 *surrogate pattern search and collocation*, Journal of Computational Physics, 229 (2010), pp. 4664–4682, doi:10.1016/
646 j.jcp.2010.03.005, <http://linkinghub.elsevier.com/retrieve/pii/S0021999110001154>.
- 647 [49] A. SHAPIRO, D. DENTCHEVA, AND A. RUSZCZYŃSKI, *Lectures on Stochastic Programming*, MPS-SIAM Series on Optimiza-
648 tion, SIAM, Philadelphia, 2009.
- 649 [50] W. F. SHARPE, *The Sharpe ratio*, J. Portfolio Management, 21 (1994), pp. 49–58.
- 650 [51] R. J. SMITH AND R. W. BLUNDELL, *An exogeneity test for a simultaneous equation Tobit model with an application to*
651 *labor supply*, Econometrica, 54 (1986), pp. 679–685.
- 652 [52] T. A. SRIVER, J. W. CHRISSIS, AND M. A. ABRAMSON, *Pattern search ranking and selection algorithms for mixed variable*
653 *simulation-based optimization*, European J. Oper. Res., 198 (2009), pp. 878–890.
- 654 [53] L. TAYLOR AND T. OTSU, *Estimation of nonseparable models with censored dependent variables and endogenous regressors*,
655 *Econom. Rev.*, (to appear).
- 656 [54] J. L. TEUGELS, *The class of subexponential distributions*, Ann. Probab., (1975), pp. 1000–1011.
- 657 [55] A. TOTH, J. A. ELLIS, T. EVANS, S. HAMILTON, C. T. KELLEY, R. PAWLOWSKI, AND S. SLATTERY, *Local improvement*
658 *results for Anderson acceleration with inaccurate function evaluations*, SIAM J. Sci. Comp., 39 (2017), pp. S47–S65,
659 doi:10.1137/16M1080677.
- 660 [56] R. VERSHYNIN, *Introduction to the non-asymptotic analysis of random matrices*, in Compressed Sensing: Theory and
661 Applications, Y. C. Eldar and G. Kutyniok, eds., Cambridge University Press, Cambridge, 2012, pp. 210–268.
- 662 [57] M. WAINWRIGHT, *High-dimensional statistics: A non-asymptotic viewpoint*, 2017, [https://www.stat.berkeley.edu/
663 ~mjwain/stat210b/Chap2_TailBounds_Jan22_2015.pdf](https://www.stat.berkeley.edu/~mjwain/stat210b/Chap2_TailBounds_Jan22_2015.pdf).
- 664 [58] J. WILLERT, X. CHEN, AND C. T. KELLEY, *Newton’s method for Monte Carlo-based residuals*, SIAM J. Numer. Anal., 53
665 (2015), pp. 1738–1757, doi:10.1137/130905691.
- 666 [59] B. YU, *Rates of convergence for empirical processes of stationary mixing sequences*, Ann. Probab., 22 (1994), pp. 94–116.