



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

## Improved algorithms for convex minimization in relative scale

**Citation for published version:**

Richtarik, P 2011, 'Improved algorithms for convex minimization in relative scale', *Siam journal on optimization*, vol. 21, no. 3, pp. 1141-1167. <https://doi.org/10.1137/090747142>

**Digital Object Identifier (DOI):**

[10.1137/090747142](https://doi.org/10.1137/090747142)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Publisher's PDF, also known as Version of record

**Published In:**

Siam journal on optimization

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



## IMPROVED ALGORITHMS FOR CONVEX MINIMIZATION IN RELATIVE SCALE\*

PETER RICHTÁRIK†

**Abstract.** In this paper we propose two modifications to Nesterov’s algorithms for minimizing convex functions in relative scale. The first is based on a bisection technique and leads to improved theoretical iteration complexity, and the second is a heuristic for avoiding restarting behavior. The fastest of our algorithms produces a solution within relative error  $O(1/k)$  of the optimum, with  $k$  being the iteration counter.

**Key words.** convex optimization, relative scale, sublinearity, Nesterov’s smoothing technique, Löwner–John ellipsoids

**AMS subject classifications.** 62K05, 65K05, 68Q25, 90C06, 90C25, 90C47, 90C60

**DOI.** 10.1137/090747142

**1. Introduction.** The theory of modern convex optimization almost uniformly assumes boundedness of the feasible set. This assumption is usually artificially enforced even for naturally unconstrained problems via the so-called “big  $M$ ” method. A clear advantage of dealing with bounded sets is the availability of a *scale* in which one can measure the absolute accuracy of a solution. However, care is needed when choosing the size of the artificially imposed bounds: large feasible sets tend to slow algorithms down, whereas small sets may lead to the exclusion of minimizers. Since there is no natural absolute scale for measuring the solutions of an unconstrained problem, it seems to be reasonable to be looking for solutions that are approximately optimal in *relative scale*. Although results of this type are rare in the convex optimization literature, some work has recently been done in this area [15], [16], [18], [19]. This contrasts with the enormous literature on combinatorial optimization where approximation algorithms are studied extensively.

In particular, Nesterov [15] showed that the above obstacles can be overcome when minimizing convex *homogeneous* functions over an affine subspace. The essence of his approach involves computing an *ellipsoidal rounding* of the subdifferential of the objective function at the origin. This family of problems encompasses essentially all unconstrained convex minimization problems via a dimension-lifting procedure. However, certain assumptions about the ellipsoidal rounding effectively limit the class of problems that can be treated.

**1.1. Contribution.** In this paper we improve the algorithms of Nesterov [14], [15] for solving unconstrained convex minimization problems within a prescribed error  $\delta$  in relative scale. We propose two modifications of the original method: the first is based on a bisection technique and leads to improved theoretical iteration complexity. The second is a heuristic for avoiding certain restarting behavior of the method. The fastest of our

---

\*Received by the editors January 20, 2009; accepted for publication (in revised form) June 28, 2011; published electronically September 29, 2011. The results of this paper were obtained in the years 2005 and 2006 while the author was a research assistant at Cornell University, working under the guidance of Mike Todd. They form a part of the PhD thesis [17, Chapter 2] of the author and were not previously published. This research was partially supported by the NSF through grants DMS-0209457 and DMS-0513337 and by the ONR through grant N00014-02-0057.

<http://www.siam.org/journals/siopt/21-3/74714.html>

†Center for Operations Research and Econometrics (CORE) and Department of Mathematical Engineering (INMA), Université catholique de Louvain, B-1348 Louvain-la-Neuve, Belgium (peter.richtarik@uclouvain.be).

algorithms produces a solution within relative error  $O(1/k)$  of the optimum, with  $k$  being the iteration counter. The bisection idea was independently used by Chudak and Eleutério [5] for obtaining the same theoretical improvement in complexity in the context of several combinatorial problems.

**1.2. Contents.** The paper is organized as follows. In section 2 we formulate the central sublinear minimization problem and briefly describe a dimension-lifting procedure for converting an unconstrained minimization instance into a linearly constrained sublinear minimization instance. Section 3 is devoted to defining basic notions and deriving key consequences of the necessary preprocessing stage of our algorithms: the computation of a pair of Löwner–John ellipsoids of a certain set. The next two parts are devoted to the description and analysis of algorithms. In section 4 we describe methods based on a simple subgradient subroutine. We first summarize Nesterov’s results and then improve them by incorporating a bisection speedup idea. We also modify the methods, at no or only negligible cost in the theoretical complexity, to allow for a “nonrestarting” behavior. In section 5 we propose more efficient methods, which are grounded in Nesterov’s smoothing technique. These are of an order of magnitude faster than those based on the subgradient routine. Next follows section 6 in which we briefly summarize the theoretical complexities and remark on the scaling invariance of the methods. In section 7 we describe several special cases to which the methods of this paper apply. The final section is devoted to computational experiments.

**1.3. Notation.** Throughout the paper,  $\mathbf{E}$  (possibly with subscripts) is a finite-dimensional real vector space and  $\mathbf{E}^*$  is its dual, i.e., the space of all linear functionals on  $\mathbf{E}$ . The action of  $g \in \mathbf{E}^*$  on  $x \in \mathbf{E}$  is written as  $\langle g, x \rangle$ . Coordinates of a vector  $y \in \mathbf{R}^l$  are denoted by superscripts in brackets; for example,  $y = (y^{(1)}, \dots, y^{(l)})$ , whereas subscripts designate vector labels. By  $\mathbf{R}_+^l$  we mean the nonnegative orthant of  $\mathbf{R}^l$ . More notation is introduced at the relevant spot in the text.

**2. Sublinear minimization.** The central problem of this paper is

$$(P) \quad \varphi^* \stackrel{\text{def}}{=} \min_{x \in \mathcal{L}} \varphi(x),$$

where  $\mathcal{L}$  is an affine subspace of a finite-dimensional real vector space  $\mathbf{E}$  not containing the origin and  $\varphi: \mathbf{E} \rightarrow \mathbf{R}$  is a *sublinear* function: convex and (positively) homogeneous of degree one. The last property means that the function is linear on every ray emanating from the origin:  $\varphi(\tau x) = \tau \varphi(x)$  for all  $\tau \geq 0$  and  $x \in \mathbf{E}$ . Note that convexity and homogeneity imply subadditivity. Define  $n := \dim \mathbf{E} = \dim \mathbf{E}^*$ .

We will further make the assumption that the zero vector lies in the interior of the (convex) subdifferential of  $\varphi$  evaluated at the origin:

$$(2.1) \quad 0 \in \text{int } \partial\varphi(0).$$

Given the properties of  $\varphi$ , condition (2.1) essentially amounts to requiring that the origin is the *unique* global minimizer of  $\varphi$ . The above assumptions imply that  $\partial\varphi(0)$  is a full-dimensional compact and convex subset of  $\mathbf{E}^*$  and that we can write<sup>1</sup>

<sup>1</sup>There is a one-to-one correspondence between finite sublinear functions and nonempty compact convex sets via the relation  $\varphi(x) = \max\{\langle g, x \rangle : g \in \mathcal{G}\}$  (this is the *support function* of  $\mathcal{G}$ ). It then follows from the definition of the subdifferential that  $\mathcal{G} = \partial\varphi(0)$ . We refer the reader Rockafellar [20]. A detailed account of the properties of sublinear functions and subdifferentials of convex functions can be found in Chapters IV and V of Hiriart-Urruty and Lemaréchal [7]. For a more compact and up-to-date treatment see Borwein and Lewis [4, Corollary 4.2.3].

$$(2.2) \quad \varphi(x) = \max\{\langle g, x \rangle : g \in \partial\varphi(0)\}.$$

That is,  $\varphi$  is the *support function* of its subdifferential at the origin. For geometric understanding of the situation implied by the assumptions, it is helpful to note that the epigraph of  $\varphi$  is a convex cone in  $\mathbf{E} \times \mathbf{R}_+$ , whose only intersection with  $\mathbf{E} \times \{0\}$  is the origin.

**2.1. Approximate solutions.** Our aim is to find an approximate solution of (P), within relative error  $\delta$ . The formal definition of the concept follows.

DEFINITION 1. A point  $x \in \mathcal{L}$  is a  $\delta$ -approximate solution to (P) if

$$\varphi(x) \leq (1 + \delta)\varphi^*.$$

In proving theorems we will often use the equivalent inequality  $\varphi(x) - \varphi^* \leq \frac{\delta}{1+\delta}\varphi(x)$ .

**2.2. Treating unconstrained convex minimization.** The general unconstrained convex minimization problem can be reformulated as a constrained sublinear problem. Let us briefly describe the construction. If  $\phi : \mathbf{E} \rightarrow \mathbf{R}$  is a convex function, its *perspective* is the function  $\varphi : \mathbf{E} \times \mathbf{R}_{++} \rightarrow \mathbf{R}$  defined by

$$\varphi(x) \stackrel{\text{def}}{=} \varphi(y, \tau) = \tau\phi(y/\tau).$$

The function  $\varphi$  is clearly linear on every feasible ray leaving from the origin. In fact, it can be shown that  $\varphi$  is convex on its domain (Hiriart-Urruty and Lemaréchal [7, Proposition 2.2.1]). It is not in general possible to extend  $\varphi$  onto the entire space  $\mathbf{E} \times \mathbf{R}$  if we want to preserve both convexity and finiteness. However, there are at least some important classes of functions for which this can be done. Consider the following example.

*Example 1.* Define  $\phi(y) = \max\{|\langle a_i, y \rangle + b^{(i)}| : i = 1, 2, \dots, m\}$ , where  $y \in \mathbf{E}$ ,  $a_1, \dots, a_m \in \mathbf{E}^*$ , and  $b \in \mathbf{R}^m$ . If we let  $x = (y, \tau)$  and  $a'_i = (a_i, b^{(i)})$  for  $i = 1, 2, \dots, m$ , then for  $\tau > 0$  we obtain

$$\varphi(x) = \tau\phi(y/\tau) = \tau \max_{1 \leq i \leq m} |\langle a_i, y/\tau \rangle + b^{(i)}| = \max_{1 \leq i \leq m} |\langle a_i, y \rangle + b^{(i)}\tau| = \max_{1 \leq i \leq m} |\langle a'_i, x \rangle|,$$

where the last equality defines a new inner product on  $\mathbf{E} \times \mathbf{R}$ . Clearly,  $\varphi$  can be extended to a sublinear function defined on the entire space. Assumption (2.1) will be satisfied if  $0 \in \text{int } \partial\varphi(0) = \text{conv}\{\pm a'_i : i = 1, 2, \dots, m\}$ .

**3. Ellipsoidal rounding and key inequalities.** As a preprocessing phase, Nesterov [15] first finds a positive definite operator  $G : \mathbf{E} \rightarrow \mathbf{E}^*$  giving rise to a pair of central ellipsoids in  $\mathbf{E}^*$ , one being contained in  $\partial\varphi(0)$  and the other containing it. This can be done using Khachiyan’s algorithm [9] or the recent method of Ahipasaoglu, Sun, and Todd [1]. We thus assume that  $G$  and  $\rho \geq 1$  are available such that

$$(3.1) \quad \mathcal{B}(G, 1) \subseteq \partial\varphi(0) \subseteq \mathcal{B}(G, \rho),$$

where  $\mathcal{B}(G, \gamma) \stackrel{\text{def}}{=} \{g \in \mathbf{E}^* : \sqrt{\langle g, G^{-1}g \rangle} \leq \gamma\}$  defines an ellipsoid in  $\mathbf{E}^*$  of radius  $\gamma$ .

The iteration complexities of the algorithms of this paper depend on the parameter  $\rho$  characterizing the quality of the ellipsoidal rounding (3.1). The following result, a celebrated theorem of John [8], gives lower bounds on the quality of rounding admitted by full-dimensional convex sets.

PROPOSITION 2 (John [8]). *Any convex body  $Q \subset \mathbf{E}^*$  admits a rounding by concentric ellipsoids with  $\rho \leq \dim \mathbf{E}^*$ . If  $Q$  is centrally symmetric, then there exists a rounding with  $\rho \leq \sqrt{\dim \mathbf{E}^*}$ .*

To see that the above result gives tight bounds, consider the following simple example.

*Example 2.* The rounding obtained by the inscribed and circumscribed balls of

- (i) a regular  $n$ -simplex has quality  $\rho = n$ ,
- (ii) the  $n$ -cube has quality  $\rho = \sqrt{n}$ .

For recent work related to ellipsoidal rounding see Belloni and Freund [2] and the references therein.

**3.1. Geometry induced by rounding.** The rounding operator  $G$  defines an inner product on  $\mathbf{E}$  via  $\langle x, y \rangle_G := \langle Gx, y \rangle$ , which in turn induces the norm  $\|x\|_G := \sqrt{\langle x, x \rangle_G}$ . The dual space  $\mathbf{E}^*$  can be equipped with the dual norm  $\|g\|_G^* := \sqrt{\langle g, G^{-1}g \rangle}$ . Notice that these norms are themselves sublinear functions and as such admit a representation similar to (2.2):

$$(3.2) \quad \|x\|_G = \max\{\langle g, x \rangle : \|g\|_G^* \leq 1\}$$

with  $\partial\|\cdot\|_G(0) = \{g \in \mathbf{E}^* : \|g\|_G^* \leq 1\}$ , and

$$(3.3) \quad \|g\|_G^* = \max\{\langle g, x \rangle : \|x\|_G \leq 1\}$$

with  $\partial\|\cdot\|_G^*(0) = \{x \in \mathbf{E} : \|x\|_G \leq 1\}$ . Also observe that the first and last sets in (3.1) are balls in  $\mathbf{E}^*$ , with respect to the dual norm, of radii 1 and  $\rho$ , respectively.

**3.2. Subgradients in the primal space.** By defining

$$\partial_G \varphi(x) \stackrel{\text{def}}{=} \{h \in \mathbf{E} : \varphi(y) \geq \varphi(x) + \langle h, x \rangle_G \text{ for all } y \in \mathbf{E}\},$$

the subgradients of  $\varphi$  can be thought of as being elements of  $\mathbf{E}$  as opposed to elements of  $\mathbf{E}^*$ . This will enable us to talk about taking steps in  $\mathbf{E}$  in the “direction” of a negative subgradient. There is a one-to-one correspondence linking the two concepts:

$$(3.4) \quad \partial_G \varphi(x) = G^{-1}[\partial\varphi(x)].$$

**3.3. Inequalities.** In view of (2.2) and (3.2), taking the maximum of the linear functional  $\langle \cdot, x \rangle$  over the sets in (3.1) gives

$$(3.5) \quad \|x\|_G \leq \varphi(x) \leq \rho \|x\|_G, \quad x \in \mathbf{E},$$

which together with subadditivity of  $\varphi$  implies that  $\varphi$  is  $\rho$ -Lipschitz:

$$\varphi(x+h) \leq \varphi(x) + \varphi(h) \leq \varphi(x) + \rho \|h\|_G.$$

From now on we will denote by  $x^*$  an arbitrary optimal solution of  $(P)$  and by  $x_0$  the minimum norm element of the feasible region—the projection of the origin onto  $\mathcal{L}$ . From (3.5) we then obtain

$$(3.6) \quad \frac{\varphi(x_0)}{\rho} \leq \|x_0\|_G \leq \|x^*\|_G \leq \varphi^* \leq \varphi(x_0) \leq \rho \|x_0\|_G.$$

Since  $\|x^* - x_0\|_G = \sqrt{\|x^*\|_G^2 - \|x_0\|_G^2}$  and  $x_0 \neq 0$  due to the assumption that  $\mathcal{L}$  does not pass through the origin, we also obtain

$$(3.7) \quad \|x^* - x_0\|_G < \|x^*\|_G \leq \varphi^* \leq \varphi(x), \quad x \in \mathcal{L}.$$

**4. Algorithms based on a subgradient subroutine.** Subgradient algorithms were studied intensively in the 1960s and 1970s by a number of researchers, among them Ermoliev, Polyak, and Shor. For comprehensive texts we refer the reader to Shor [21] and Goffin [6]. For our purposes we will manage with a result about the performance of a standard constant step-length subgradient algorithm applied to a convex Lipschitz function [12, section 3.2.3].

**4.1. A constant step-length subgradient algorithm.** The subgradient algorithm we are going to describe works in a more general setting than that of problem  $(P)$ . For the sake of this subsection only, consider the problem of minimizing a convex Lipschitz continuous function  $\varphi : \mathbf{E} \rightarrow \mathbf{R}$  with Lipschitz constant  $\gamma$  over a *simple* closed convex set  $Q_1$ :

$$(P_{sg}) \quad \varphi^* \stackrel{\text{def}}{=} \min\{\varphi(x) : x \in Q_1\}.$$

By simple set we mean one allowing for easy computation of projections onto it (symbol  $\text{proj}$  will denote the projection operator). In this setting  $\mathbf{E}$  is assumed to be equipped with an inner product. Problem  $(P)$  is a special case of  $(P_{sg})$  with

- $\varphi$  having additional properties,
- $\gamma = \rho$  and  $Q_1 = \mathcal{L}$ , and
- $\mathbf{E}$  made Euclidean by the introduction of the inner product induced by  $G$ .

The following is a standard result (see, for example, Nesterov [12, Theorem 3.2.2]).

**PROPOSITION 3.** *If  $\|x^* - x_0\| \leq R$  for some  $x_0 \in \mathbf{E}$ , minimizer  $x^*$  of  $(P_{sg})$ , and  $R > 0$ , then the output  $x = \text{Subgrad}(\varphi, Q_1, x_0, R, N)$  of Algorithm 1 run on an instance of problem  $(P_{sg})$  satisfies*

$$(4.1) \quad \varphi(x) - \varphi^* \leq \frac{\gamma R}{\sqrt{N+1}}.$$

---

**ALGORITHM 1. (SUBGRAD) CONSTANT STEP-LENGTH SUBGRADIENT SCHEME.**

- 1: **Input:**  $\varphi, Q_1, x_0, R, N$ ;
  - 2:  $\kappa = \frac{R}{\sqrt{N+1}}$ ;
  - 3: **for**  $k = 0$  **to**  $N - 1$
  - 4:   pick  $g \in \partial\varphi(x_k)$  **if**  $g = 0$  **then**  $x_k$  is optimal and **exit**;
  - 5:    $x_{k+1} = \text{proj}_{Q_1}(x_k - \kappa \frac{g}{\|g\|})$ ;
  - 6: **end for**
  - 7: **Output:**  $x_k$  with best objective value
- 

For Proposition 3 to hold it suffices to require that  $\varphi$  be Lipschitz on the ball around  $x^*$  with radius  $R$ .

**4.2. Basic algorithmic ideas.** As the previous subsection indicates, the basic idea for solving  $(P)$  will be that of using the subgradient method (Algorithm 1). The main issue with this algorithm, apart from the fact that it is slow (it requires  $O(1/\epsilon^2)$

iterations to output an  $\epsilon$ -optimal solution in the additive sense), is the need to supply an initial point  $x_0$  and an upper bound  $R$  on  $\|x^* - x_0\|$ .

The particular choice of  $x_0$  as the projection of the origin onto the feasible set of  $(P)$  makes sense for at least two reasons. First, notice that if the ellipsoidal rounding of  $\partial\varphi(0)$  is perfectly tight ( $\rho = 1$ ), then by (3.5) we have  $\varphi(x) \equiv \|x\|_G$ , and therefore  $x_0$  is the optimal solution of  $(P)$ . In fact, notice that by (3.6),

$$(4.2) \quad \varphi(x_0) \leq \rho\varphi^*,$$

and hence  $x_0$  is a  $(\rho - 1)$ -approximate solution of  $(P)$ . The better the rounding, the better the approximation factor. Second, (3.7) offers the readily available upper bound  $R = \varphi(x_0)$ . Of course,  $\varphi^*$  would be better; the issue is that it is not known.

**4.2.1. Good but unavailable upper bound.** Let us formally apply Algorithm 1 to  $(P)$  with  $R = \varphi^*$ . To achieve the required relative accuracy, it then suffices to run it for  $N = \lfloor \rho^2 / \delta^2 \rfloor$  iterations because, by Proposition 3,

$$\varphi(x) - \varphi^* \leq \frac{\rho R}{\sqrt{N+1}} \leq \frac{\rho\varphi^*}{\sqrt{\rho^2/\delta^2}} = \delta\varphi^*.$$

**4.2.2. Available but bad upper bound.** Since the previous upper bound is unknown, it seems reasonable to instead use the worse (but available) bound  $R = \varphi(x_0)$ . If we wish to guarantee a solution within relative error  $\delta$ , we need to take  $N = \lfloor \rho^4 / \delta^2 \rfloor$  iterations. The argument is exactly the same and uses (4.2).

**4.2.3. Iteratively updated upper bound.** To move toward the better of the two extremes, Nesterov [15] proposed a scheme (Algorithm 2) which uses the subgradient method as a subroutine, iteratively decreasing the known upper bound. This algorithm starts by running the subgradient method for  $O(\rho^2 / \delta^2)$  iterations with the available upper bound  $\varphi(x_0)$ . In the case when the subgradient subroutine is doing well and manages to decrease the objective value by a constant factor, the previously available upper bound also decreases by the same factor. This improved bound is then used to run the next subgradient subroutine, again starting from  $x_0$ .

---

**ALGORITHM 2. (SUBSEARCH) SUBGRADIENT SEARCH SCHEME.**

- 1: **Input:**  $\varphi, \mathcal{L}, x_0, \rho, \beta > 0, \delta$ ;
  - 2:  $\hat{x}_0 = x_0, c = e^\beta, k = 1$ ;
  - 3:  $N = \lfloor c^2 \rho^2 (1 + \frac{1}{\delta})^2 \rfloor$ ;
  - 4:  $\hat{x}_k = \text{Subgrad}(\varphi, \mathcal{L}, x_0, \varphi(\hat{x}_{k-1}), N)$ ;
  - 5: **while**  $\varphi(\hat{x}_k) < \varphi(\hat{x}_{k-1})/c$  **do**
  - 6:      $k = k + 1$ ;
  - 7:      $\hat{x}_k = \text{Subgrad}(\varphi, \mathcal{L}, x_0, \varphi(\hat{x}_{k-1}), N)$ ;
  - 8: **end while**
  - 9: **Output:**  $\hat{x}_k$
- 

The performance of Algorithm 2 is substantially better than the naive one-time application of the subgradient method with the bad but available upper bound. However, it underperforms the one-time application of the subgradient method with the good but unknown upper bound by a factor of  $O(\ln \rho)$ . The performance of the method, as

analyzed by Nesterov [15], is summarized in Proposition 4. We include the proof because it is short and offers insight into the subsequent improvements we propose in the following subsections. We will also refer to parts of it later.

PROPOSITION 4 (Nesterov [15, Theorem 3]). *Algorithm 2 returns a  $\delta$ -approximate solution of (P) and takes at most*

$$(4.3) \quad e^{2\beta} \rho^2 \left(1 + \frac{1}{\delta}\right)^2 \left(1 + \frac{1}{\beta} \ln \rho\right)$$

steps of the subgradient method. If  $\beta$  is a constant, then the number of steps is

$$(4.4) \quad O\left(\frac{\rho^2}{\delta^2} \ln \rho\right).$$

The optimal choice is  $\beta = \frac{1}{2}(\sqrt{t^2 + 2t} - t) \approx \frac{1}{2}$ , with  $t = \ln \rho$ .

*Proof.* Assume that the algorithm stops at iteration  $k$ , failing to satisfy the “while” clause at step 5. In view of (3.6) we have

$$\frac{\varphi(x_0)}{\rho} \leq \varphi^* \leq \varphi(\hat{x}_{k-1}) < \frac{\varphi(x_0)}{e^{\beta(k-1)}},$$

and by comparing the first and the last term in this chain of inequalities we conclude that the number of calls of the subgradient subroutine is at most  $1 + \beta^{-1} \ln \rho$ . The bound (4.3) is obtained by multiplying this by  $N$  from step 3 of the algorithm. Minimizing (4.3) in  $\beta$  gives the final statement. It remains to show that the output is as specified. Indeed, using the termination rule from step 5 and applying Proposition 3 to the last call of the subgradient subroutine, we get

$$\varphi(\hat{x}_k) - \varphi^* \leq \frac{\rho \varphi(\hat{x}_{k-1})}{\sqrt{N+1}} \leq \frac{\rho e^\beta \varphi(\hat{x}_k)}{\sqrt{N+1}} \leq \frac{\delta}{1+\delta} \varphi(\hat{x}_k) \quad \square$$

**4.3. Bisection improvement.** Each outer iteration of Algorithm 2, possibly except the last one, produces a *guaranteed* upper bound on the distance of  $x_0$  from the set of minimizers of (P)—better by a constant factor than the one available before. Loosely speaking, we will show that by allowing for *guesswork* it is possible to improve the theoretical performance of this algorithm (the same improvement was independently obtained by Chudak and Eleutério [5] in the context of combinatorial applications). The key observation is formulated in the following lemma.

LEMMA 5. *If  $\varphi^* \leq R$  and  $N = \lfloor \rho^2 / \beta^2 \rfloor$  for some  $\beta > 0$ , then*

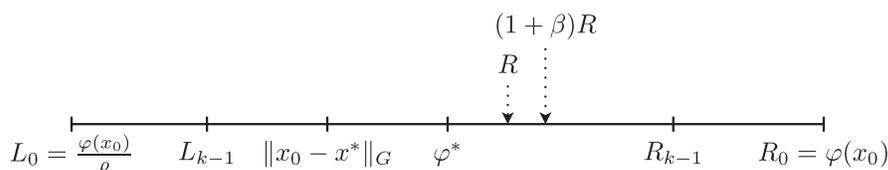
$$x = \text{Subgrad}(\varphi, \mathcal{L}, x_0, R, N)$$

satisfies

$$(4.5) \quad \varphi(x) - \beta R \leq \varphi^* \quad \text{and} \quad \varphi(x) \leq (1 + \beta)R.$$

*Proof.* By Proposition 3 we have  $\varphi(x) - \varphi^* \leq \rho R / \sqrt{N+1} \leq \beta R$ , and hence  $\varphi(x) \leq \varphi^* + \beta R \leq (1 + \beta)R$ .  $\square$

The above result essentially states that for *any* positive  $R$  we can, at the cost of  $O(\rho^2 / \beta^2)$  iterations of the subgradient method (Algorithm 1), either get a certificate

FIG. 1. Bisection step  $k$ .

that  $\varphi^* \leq (1 + \beta)R$  (if  $x$  satisfies  $\varphi(x) \leq (1 + \beta)R$ ) or that  $R < \varphi^*$  (if  $\varphi(x) > (1 + \beta)R$ ). In any case, we either get an upper or lower bound on  $\varphi^*$ .

Note that thanks to (3.6), we are in the possession of an initial lower and upper bound on  $\varphi^*$ : if we set  $L_0 = \|x_0\|_G$  and  $R_0 = \varphi(x_0)$ , then  $\frac{\varphi(x_0)}{\rho} \leq L_0 \leq \varphi^* \leq R_0$ , with

$$(4.6) \quad \frac{R_0}{L_0} \leq \rho.$$

Assume that at step  $k$  we have  $L_k \leq \varphi^* \leq R_k$ , with  $q_k \stackrel{\text{def}}{=} R_k/L_k > 1 + \beta$  (see Figure 1). Pick  $R$  so that

$$(4.7) \quad L_k < R < (1 + \beta)R < R_k.$$

For this  $R$  let  $x$  be given by Lemma 5. There are two possibilities. If  $\varphi(x) \leq (1 + \beta)R$ , then in view of (4.5) we can update

$$(4.8) \quad L_{k+1} = \max\{\varphi(x) - \beta R, L_k\}, \quad R_{k+1} = \min\{R_k, \varphi(x)\} \leq (1 + \beta)R.$$

If  $\varphi(x) > (1 + \beta)R$ , then we can set

$$(4.9) \quad L_{k+1} = R, \quad R_{k+1} = \min\{R_k, \varphi(x)\}.$$

This bisection procedure is then repeated until  $q_k < (1 + \tau)(1 + \beta)$  for some  $\tau > 0$ . The following lemma states how much improvement in  $q_k$  can be obtained by a single bisection step.

LEMMA 6. Assume  $L_k \leq \varphi^* \leq R_k$ ,  $q_k > 1 + \beta$ , and  $\beta > 0$ . After a single bisection step with  $R = [L_k R_k / (1 + \beta)]^{1/2}$ , we obtain  $L_{k+1} \leq \varphi^* \leq R_{k+1}$  satisfying

$$(4.10) \quad q_{k+1} \leq (1 + \beta)^{1/2} q_k^{1/2}.$$

*Proof.* It is easy to see that (4.7) holds. Observing that  $R$  is chosen so that  $(1 + \beta)R/L_k = R_k/R$ , in view of (4.8) and (4.9) we obtain

$$q_{k+1} = \frac{R_{k+1}}{L_{k+1}} \leq \max \left\{ \frac{(1 + \beta)R}{L_k}, \frac{\min\{R_k, \varphi(x)\}}{R} \right\} \leq \frac{R_k}{R} = (1 + \beta)^{1/2} q_k^{1/2}. \quad \square$$

The ideas outlined above lead to Algorithm 3, whose performance is analyzed in the next theorem.

---

**ALGORITHM 3. (SubBis) SUBGRADIENT BISECTION SCHEME.**

- 1: **Input:**  $\varphi, \mathcal{L}, x_0, \rho, \beta, \tau, \delta$ ;
  - 2:  $k = 0, L_0 = \|x_0\|_G, R_0 = \varphi(x_0), c = (1 + \tau)(1 + \beta), N = \lfloor \rho^2 / \beta^2 \rfloor$ ;
  - 3: **while**  $R_k / L_k > c$  **do**
  - 4:  $R = \sqrt{\frac{L_k R_k}{1 + \beta}}, x = \text{Subgrad}(\varphi, \mathcal{L}, x_0, R, N)$ ;
  - 5: **if**  $\varphi(x) \leq (1 + \beta)R$  **then**
  - 6: set  $R_{k+1}, L_{k+1}$  as in (4.8)
  - 7: **else**
  - 8: set  $R_{k+1}, L_{k+1}$  as in (4.9)
  - 9: **end if**
  - 10:  $k = k + 1$
  - 11: **end while**
  - 12:  $N = \lfloor \frac{R_k^2}{L_k^2} \rho^2 (1 + \frac{1}{\delta})^2 \rfloor, x = \text{Subgrad}(\varphi, \mathcal{L}, x_0, R_k, N)$ ;
  - 13: **Output:**  $x$
- 

THEOREM 7. *Algorithm 3 returns a  $\delta$ -approximate solution of (P) and takes at most*

$$(4.11) \quad \frac{\rho^2}{\beta^2} \left[ 1 + \log_2 \left( \frac{\ln \rho}{\ln(1 + \tau)} \right) \right] + (1 + \tau)^2 (1 + \beta)^2 \rho^2 \left( 1 + \frac{1}{\delta} \right)^2$$

*steps of the subgradient subroutine. If  $\beta$  is a constant, the number of steps is*

$$(4.12) \quad O \left( \rho^2 \left( \frac{1}{\delta^2} + \ln \ln \rho \right) \right).$$

*Proof.* Let us first analyze the bisection phase (the “while” loop). The repeated use of Lemma 6 gives

$$q_k \leq (1 + \beta)^{\frac{1}{2}} q_{k-1}^{\frac{1}{2}} \leq (1 + \beta)^{\frac{1}{2}} (1 + \beta)^{\frac{1}{4}} \cdots (1 + \beta)^{\frac{1}{2^k}} q_0^{\frac{1}{2^k}} \stackrel{(4.6)}{\leq} (1 + \beta) \rho^{\frac{1}{2^k}}.$$

The smallest integer  $k$  for which  $(1 + \beta) \rho^{\frac{1}{2^k}} \leq (1 + \tau)(1 + \beta)$  is  $k^* = \lceil \log_2(\ln \rho / \ln(1 + \tau)) \rceil$ , and hence the total number of lower-level subgradient method iterations of the bisection phase is at most  $N_{\text{bis}} = \rho^2 \beta^{-2} k^*$ . The statement then follows by adding  $N_{\text{bis}}$  and the number of iterations needed for the finalization phase (step 12). It remains to show that the output of the algorithm is as specified. Notice that  $L_k \leq \varphi(x)$ , and Proposition 3 applied to the subgradient method call at step 12 give

$$\varphi(x) - \varphi^* \leq \frac{\rho R_k}{\sqrt{N + 1}} \leq \frac{\rho \frac{R_k}{L_k} \varphi(x)}{\sqrt{N + 1}} \leq \frac{\delta}{1 + \delta} \varphi(x). \quad \square$$

**4.4. Nonrestarting algorithms.** Algorithms **SubSearch** and **SubBis** (Algorithms 2 and 3) use the subgradient subroutine *always started from one point*, denoted  $x_0$ , which is defined as the projection of the origin onto the feasible set. This point is indeed special, as it allows for the key inequalities (3.6) and (3.7), which in turn drive the analysis in both algorithms. The first of these inequalities makes  $x_0$  indispensable as the starting point of the very *first* subgradient subroutine call in both algorithms, making it possible to construct initial lower and upper bounds on  $\varphi^*$ . It is hard to think of a different readily computable point that could serve the same purpose.

The issue we are going to touch upon in this subsection concerns the use of  $x_0$  as the starting point in all *subsequent* calls of the subroutine. In our view, *restarting* from this particular point seems to be convenient for the sake of the proofs rather than efficient algorithmically. Let us elaborate on this a bit. Both algorithms mentioned above can be viewed as simultaneously optimizing (solving  $(P)$ ) and searching for a good upper bound on  $\|x_0 - x^*\|_G$  in order to look less like the “do-it-all-with-the-available-but-bad-upper-bound” and more like the “do-it-all-with-the-good-but-unavailable-upper-bound” algorithm. Combining these two goals is possible because  $\varphi^*$  is both the optimal value of  $(P)$  and an upper bound on  $\|x_0 - x^*\|_G$ . It seems likely that the optimization goal could be attained faster if we could use the current best point, as opposed to  $x_0$ , to start every call of the subroutine. Although both algorithms gather information about increasingly better iterates  $\{\hat{x}_k\}$ , this knowledge is used only to update the upper bound on  $\|x_0 - x^*\|_G$  in the next call of the subgradient subroutine and not to start the subroutine itself from a better point. There is a good reason for that though: even if some point  $\hat{x}_k$  obtained along the way in one of the algorithms was much better than  $x_0$  in terms of its objective value, there are no theoretical guarantees that  $\|\hat{x}_k - x^*\|_G$  will be smaller. Starting the subgradient subroutine from such a point thus means combining a probable advantage with a possible disadvantage. A simple observation reveals that the disadvantage factor is under control. Indeed, for any  $x \in \mathcal{L}$ ,

$$(4.13) \quad \|x - x^*\|_G \leq \|x\|_G + \|x^*\|_G \stackrel{(3.7)}{\leq} \|x\|_G + \varphi^* \leq \|x\|_G + \varphi(x) \stackrel{(3.5)}{\leq} 2\varphi(x).$$

This means that whenever the subgradient method outputs some point  $x$ , we have an upper bound on  $\|x - x^*\|_G$  available. Therefore, on the next call we can run the method starting at  $x$  with  $R = \|x\|_G + \varphi(x)$ .

**4.4.1. Nonrestarting version of SubSearch.** Algorithm 4 is a modified version of Algorithm 2 in the spirit of the preceding discussion. The theoretical performance is unchanged.

---

**ALGORITHM 4. (SUBSEARCHNR) NONRESTARTING SUBGRADIENT SEARCH SCHEME.**

- 1: **Input:**  $\varphi, \mathcal{L}, x_0, \rho, \delta$ ;
  - 2:  $\hat{x}_0 = x_0, c = \sqrt{e}, k = 1$ ;
  - 3:  $N = \lfloor c^2 \rho^2 (1 + \frac{1}{\delta})^2 \rfloor, N' = \lfloor 4c^2 \rho^2 (1 + \frac{1}{\delta})^2 \rfloor, R = \varphi(\hat{x}_0)$ ;
  - 4:  $\hat{x}_k = \text{Subgrad}(\varphi, \mathcal{L}, \hat{x}_0, R, N)$ ;
  - 5: **while**  $\varphi(\hat{x}_k) < \varphi(\hat{x}_{k-1})/c$  **do**
  - 6:      $k = k + 1$ ;
  - 7:      $R = \|\hat{x}_{k-1}\|_G + \varphi(\hat{x}_{k-1})$ ;
  - 8:      $\hat{x}_k = \text{Subgrad}(\varphi, \mathcal{L}, \hat{x}_{k-1}, R, N')$ ;
  - 9: **end while**
  - 10: **Output:**  $\hat{x}_k$
- 

**THEOREM 8.** *Algorithm 4 outputs a  $\delta$ -approximate solution of  $(P)$ . The number of calls of the subgradient subroutine is at most  $1 + 2 \ln \rho$ , and the total number of lower-level subgradient steps is hence at most*

$$(4.14) \quad 4e\rho^2 \left(1 + \frac{1}{\delta}\right)^2 (1 + 2 \ln \rho) = O\left(\frac{\rho^2}{\delta^2} \ln \rho\right).$$

*Proof.* The proof of the upper bound on the number of the outer level iterations is exactly the same as for Algorithm 2. If the algorithm terminates with  $k = 1$ , it is identical to Nesterov's, and the result follows (we can drop the constant 4 in this case). For  $k > 1$ , the analysis is analogous:

$$\varphi(\hat{x}_k) - \varphi^* \leq \frac{\rho R}{\sqrt{N+1}} \stackrel{(4.13)}{\leq} \frac{\rho 2\varphi(\hat{x}_{k-1})}{2c\rho\left(1 + \frac{1}{\delta}\right)} \leq \frac{\delta}{1 + \delta} \varphi(\hat{x}_k). \quad \square$$

**4.4.2. Nonrestarting bisection algorithm.** The following fact plays the role of Lemma 5 in the design and analysis of a nonrestarting bisection algorithm (Algorithm 5).

LEMMA 9. *Let  $x' \in \mathcal{L}$  and assume  $\varphi^* \leq R$ . If we let  $N = \lfloor \rho^2 / \beta^2 \rfloor$  for some  $\beta > 0$ , then*

$$x = \text{Subgrad}(\varphi, \mathcal{L}, x', R + \|x'\|_G, N)$$

satisfies

$$(4.15) \quad \varphi(x) - \beta(\|x'\|_G + R) \leq \varphi^* \quad \text{and} \quad \varphi(x) \leq (1 + \beta)R + \beta\|x'\|_G.$$

*Proof.* By (4.13) we have  $\|x' - x^*\|_G \leq \|x'\|_G + R$ , and hence by Proposition 3,

$$\varphi(x) - \varphi^* \leq \rho \frac{\|x'\|_G + R}{\sqrt{N+1}} \leq \beta(\|x'\|_G + R).$$

Rearranging the expression gives the first inequality in (4.15); the second inequality follows from (3.5).  $\square$

The idea of updating the lower and upper bounds is analogous to the restarting version of the algorithm. Assume that at step  $k$  we have  $L_k \leq \varphi^* \leq R_k$ , with  $q_k = R_k/L_k > 1 + \beta$ . Pick  $R$  so that (4.7) holds and  $x'$  such that  $R_k = \varphi(x')$ , and let  $x$  be given by Lemma 9. Again, we have two possibilities. Notice that if  $\varphi(x) \leq (1 + \beta)R + \beta\|x'\|_G$ , then since  $\|x'\|_G \leq \varphi(x') = R_k$ , we have

$$(4.16) \quad \varphi(x) \leq (1 + \beta)R + \beta R_k \stackrel{(4.7)}{\leq} R_k,$$

as long as  $\beta \leq 1$ . We can thus update

$$(4.17) \quad L_{k+1} = \max\{\varphi(x) - \beta(\|x'\|_G + R), L_k\}, \quad R_{k+1} = \varphi(x).$$

If  $\varphi(x) > (1 + \beta)R + \beta\|x'\|_G$ , we can set

$$(4.18) \quad L_{k+1} = R, \quad R_{k+1} = \min\{R_k, \varphi(x)\}.$$

The improvement in  $q_k$  after a single bisection step is given in the following result.

LEMMA 10. *Assume  $L_k \leq \varphi^* \leq R_k$ ,  $q_k > 1 + \beta$ , and  $0 < \beta \leq 1$ . After a single bisection step with  $R = \lfloor L_k R_k / (1 + \beta) \rfloor^{1/2}$ , we obtain  $L_{k+1} \leq \varphi^* \leq R_{k+1}$  satisfying*

$$(4.19) \quad q_k \leq \left(\beta + \frac{1}{\sqrt{2}}\right) q_{k-1}.$$

*Proof.* Note that (4.7) holds. In view of (4.17), (4.16), and (4.18) we get

$$\begin{aligned}
q_{k+1} = \frac{R_{k+1}}{L_{k+1}} &\leq \max \left\{ \frac{(1+\beta)R + \beta R_k}{L_k}, \frac{\min\{R_k, \varphi(x)\}}{R} \right\} \\
&\leq \frac{(1+\beta)R + \beta R_k}{L_k} \\
&= \sqrt{1+\beta} \sqrt{q_k} + \beta q_k = \sqrt{\frac{1+\beta}{q_k}} q_k + \beta q_k < \sqrt{\frac{1}{2}} q_k + \beta q_k. \quad \square
\end{aligned}$$

THEOREM 11. Algorithm 5 run with  $\beta$  satisfying  $0 < \beta < 1 - \frac{1}{\sqrt{2}}$  returns a  $\delta$ -approximate solution of (P) and takes at most

$$(4.20) \quad \frac{\rho^2}{\beta^2} \left[ 1 + \frac{\ln \rho - \ln[(1+\tau)(1+\beta)]}{-\ln(\beta + 1/\sqrt{2})} \right] + 4(1+\tau)^2(1+\beta)^2 \rho^2 \left( 1 + \frac{1}{\delta} \right)^2$$

steps of the subgradient subroutine. If  $\tau$  and  $\beta$  are chosen as constants, this becomes

$$(4.21) \quad O\left(\frac{\rho^2}{\delta^2} + \rho^2 \ln \rho\right).$$

*Proof.* Let us first analyze the bisection phase. Repeated use of Lemma 10 gives  $q_k \leq (\beta + 1/\sqrt{2})^k q_0 \leq (\beta + 1/\sqrt{2})^k \rho$ . The smallest integer  $k$  for which the last expression drops below  $c = (1+\tau)(1+\beta)$  is  $k^* = \lceil \ln(\rho/c) / \ln(1/(\beta + 1/\sqrt{2})) \rceil = O(\ln \rho)$ , and hence the total number of lower-level subgradient method iterations of the bisection phase is  $N_{\text{bis}} = \lfloor \rho^2 / \beta^2 \rfloor k^*$ . The guarantee (4.20) follows by adding  $N_{\text{bis}}$  and  $N$  from step 13. The output of the algorithm is as specified; the analysis is identical to that in Theorem 8.  $\square$

Note that the nonrestarting version of the bisection algorithm has a slightly worse complexity bound—we have lost one logarithm in (4.21) in comparison with (4.12). However, the bisection strategy separates the  $\delta$  from the logarithmic term when compared to the bound (4.14) for the **SubSearch** algorithm.

---

**ALGORITHM 5. (SUBBISNR) NONRESTARTING SUBGRADIENT BISECTION SCHEME.**

- 1: **Input:**  $\varphi, \mathcal{L}, x_0, \rho, \beta, \tau, \delta$ ;
  - 2:  $k = 0, x' = x_0, L_0 = \|x_0\|_G, R_0 = \varphi(x_0)$ ;
  - 3:  $c = (1+\tau)(1+\beta), N = \lfloor \rho^2 / \beta^2 \rfloor$ ;
  - 4: **while**  $R_k / L_k > c$  **do**
  - 5:  $R = \sqrt{\frac{L_k R_k}{1+\beta}}, x = \text{Subgrad}(\varphi, \mathcal{L}, x', \|x'\|_G + R, N)$ ;
  - 6: **if**  $\varphi(x) \leq (1+\beta)R + \beta \|x'\|_G$  **then**
  - 7: set  $L_{k+1}, R_{k+1}$  as in (4.17)
  - 8: **else**
  - 9: set  $L_{k+1}, R_{k+1}$  as in (4.18)
  - 10: **end if**
  - 11:  $x' = x, k = k + 1$
  - 12: **end while**
  - 13:  $N = \lfloor 4 \frac{R_k^2}{L_k^2} \rho^2 (1 + \frac{1}{\delta})^2 \rfloor, x = \text{Subgrad}(\varphi, \mathcal{L}, x', \|x'\|_G + R_k, N)$ ;
  - 14: **Output:**  $x$
- 

**5. Algorithms based on smoothing.** We have seen in section 4 that problem (P) allows for simple algorithms that require  $O(\delta^{-2})$  iterations of the subgradient method.

We have improved Nesterov’s subgradient search algorithm (Algorithm 2), which needs  $O(\rho^2\delta^{-2} \ln \rho)$  iterations, by incorporating a simple bisection idea and obtained Algorithm 3 with the slightly better  $O(\rho^2\delta^{-2} + \delta^{-2} \ln \ln \rho)$  complexity. That is, we have improved the dependence on the rounding parameter  $\rho$  but not on the error parameter  $\delta$ .

We start in the following subsection by briefly describing Nesterov’s smoothing technique [13] and the implied algorithm for smooth minimization of nonsmooth functions. It is not our intention to describe the approach in full generality; rather, we will adapt the results to the setting of problem  $(P)$ —the minimization of a nonnegative sublinear (convex and homogeneous) function vanishing at the origin only.

**5.1. The setting.** Nesterov [13] considers a rather general *nonsmooth* convex optimization problem and shows that it is possible to solve it in  $O(\epsilon^{-1})$  iterations of a gradient-type method if a solution within absolute error  $\epsilon$  is sought. His novel approach involves two phases. The first is a preprocessing phase in which one approximates the objective function by a smooth function with Lipschitz continuous gradient. The second phase amounts to running an optimal smooth method [11], [12] with complexity  $O(\epsilon^{-1/2})$  applied to the smooth function.

We will describe the model for sublinear functions. Consider the following more general version of problem  $(P)$ , with  $\varphi$  replaced by an arbitrary sublinear function and  $\mathcal{L}$  (or  $\mathcal{L}$  intersected with a large ball) replaced by a compact and convex subset  $Q_1$  of  $\mathbf{E}_1 := \mathbf{E}$ :

$$(P') \quad \varphi^* := \min_x \{\varphi(x) : x \in Q_1\}.$$

Notice that  $\varphi$  can be written as

$$(5.1) \quad \varphi(x) = \max_g \{\langle g, x \rangle : g \in \partial\varphi(0)\}.$$

To allow for some modeling flexibility, the purpose of which will be clear later, we will instead consider the following family of representations of the objective function:

$$(5.2) \quad \varphi(x) = \max_y \{\langle Ax, y \rangle : y \in Q_2\}.$$

Here we are introducing a new finite-dimensional real vector space  $\mathbf{E}_2$ , a linear operator  $A : \mathbf{E}_1 \rightarrow \mathbf{E}_2^*$ , and a compact and convex set  $Q_2 \subset \mathbf{E}_2$ .

DEFINITION 12. *The adjoint of  $A$  is the operator  $A^* : \mathbf{E}_2 \rightarrow \mathbf{E}_1^*$  defined via*

$$\langle Ax, y \rangle = \langle A^*y, x \rangle \quad \text{for all } x \in \mathbf{E}_1, y \in \mathbf{E}_2.$$

We assume that the spaces  $\mathbf{E}_1$  and  $\mathbf{E}_2$  are equipped with norms  $\|\cdot\|_1$  and  $\|\cdot\|_2$ , respectively<sup>2</sup>, and the dual spaces  $\mathbf{E}_1^*$  and  $\mathbf{E}_2^*$  with the corresponding dual norms

$$(5.3) \quad \|g\|_1^* := \max\{\langle g, x \rangle : \|x\|_1 \leq 1\} \quad \text{and} \quad \|h\|_2^* := \max\{\langle h, y \rangle : \|y\|_2 \leq 1\}$$

for  $g \in \mathbf{E}_1^*$  and  $h \in \mathbf{E}_2^*$ .

---

<sup>2</sup>The numbers are subscripts referring to the spaces in which the norms are defined and are not intended to suggest the use of the  $\ell_1$  and  $\ell_2$  norms.

DEFINITION 13. *The norm of  $A$  is defined by*

$$(5.4) \quad \|A\|_{1,2} := \max_{x,y} \{\langle Ax, y \rangle : \|x\|_1 = 1, \|y\|_2 = 1\}.$$

One can similarly define  $\|A^*\|_{2,1}$ .

It follows easily from the definition that

$$(5.5) \quad \|A\|_{1,2} = \max_x \{\|Ax\|_2^* : \|x\|_1 = 1\} = \|A^*\|_{2,1} = \max_y \{\|A^*y\|_1^* : \|y\|_2 = 1\}.$$

*Example 3.* Consider the function

$$\varphi_\infty(x) := \max_i \{|\langle a_i, x \rangle| : i = 1, 2, \dots, m\},$$

where  $x \in \mathbf{E}_1 = \mathbf{R}^n$ ,  $a_i \in \mathbf{E}_1^* = \mathbf{R}^n$ , and  $\langle g, x \rangle = \sum_{i=1}^n g^{(i)} x^{(i)}$ . Note that in the following three representations of  $\varphi_\infty$  the structure of the set  $Q_2$  gets simpler as the dimension of the space  $\mathbf{E}_2$  increases.

1.  $\mathbf{E}_2 = \mathbf{E}_2^* = \mathbf{R}^n$ ,  $Q_2 = \text{conv}\{\pm a_i : i = 1, 2, \dots, m\}$ , and  $A = I$ . This seems to be the most natural and straightforward representation.
2.  $\mathbf{E}_2 = \mathbf{E}_2^* = \mathbf{R}^m$ ,  $Q_2 = \{y \in \mathbf{R}^m : \sum_{i=1}^m |y^{(i)}| \leq 1\}$ , and  $A$  is the  $m \times n$  matrix with rows  $a_1, \dots, a_m$ . In this case we have

$$\varphi_\infty(x) = \max \left\{ \sum_{i=1}^m y^{(i)} \langle a_i, x \rangle : \sum_{i=1}^m |y^{(i)}| \leq 1 \right\}.$$

3.  $\mathbf{E}_2 = \mathbf{E}_2^* = \mathbf{R}^{2m}$ ,  $Q_2$  is the unit simplex in  $\mathbf{R}^{2m}$ , and  $A$  is the  $2m \times n$  matrix with rows composed of  $a_1, \dots, a_m$  and  $-a_1, \dots, -a_m$ :

$$\varphi_\infty(x) = \max \left\{ \sum_{i=1}^m (y_1^{(i)} - y_2^{(i)}) \langle a_i, x \rangle : \sum_{i=1}^m y_1^{(i)} + y_2^{(i)} = 1, y_1^{(i)}, y_2^{(i)} \geq 0 \right\}.$$

If we let  $\theta(y) \stackrel{\text{def}}{=} \min_x \{ \langle A^*y, x \rangle : x \in Q_1 \}$ , then because both  $Q_1$  and  $Q_2$  are convex and compact and  $\langle A^*y, x \rangle \equiv \langle y, Ax \rangle$  is bilinear, we can apply a standard minimax result and rewrite  $(P')$  as follows:

$$(P'') \quad \varphi^* = \theta^* \stackrel{\text{def}}{=} \max_y \{ \theta(y) : y \in Q_2 \}.$$

**5.2. Smoothing and an efficient smooth method.** In the first phase of Nesterov's approach, the objective function of  $(P')$  is approximated by a smooth convex function with Lipschitz continuous gradient. An approximation with error  $O(\epsilon)$  has gradient with Lipschitz constant of  $O(1/\epsilon)$ . The second phase consists of applying to  $(P)$  (with the objective function replaced by its smooth approximation) an efficient smooth method (Algorithm 6) requiring  $O(1/\sqrt{\epsilon})$  iterations of a gradient type. The smooth algorithm is capable of producing points  $\hat{x}$  and  $\hat{g}$  feasible to *both*  $(P')$  and  $(P'')$ , respectively, such that  $\varphi(\hat{x}) - \theta(\hat{g}) = O(1/\epsilon)$ . Because  $\varphi^* = \theta^*$ , these points are approximate optimizers in their respective problems (in the additive sense).

The smoothing approach assumes the availability of *prox-functions*  $d_1$  and  $d_2$  for the sets  $Q_1$  and  $Q_2$ , respectively. These are continuous and strongly convex nonnegative functions defined on these sets, with convexity parameters  $\sigma_1$  and  $\sigma_2$ , respectively. Let  $x_0$  be the *center* of the set  $Q_1$  (think  $Q_1 = \mathcal{L}$ ):

$$(5.6) \quad x_0 := \arg \min_x \{d_1(x) : x \in Q_1\}.$$

We assume that  $d_1$  vanishes at its center, and hence the above properties imply

$$d_1(x) \geq \frac{1}{2} \sigma_1 \|x - x_0\|_1^2.$$

For example, if  $d_1(x) := \frac{1}{2} \|x\|_1^2$  (so  $\sigma_1 = 1$ ) and  $Q_1$  is the intersection of  $\mathcal{L}$  and a large-enough ball centered at the origin, then  $x_0$  coincides with its earlier definition. Notice that for  $d_1(x) = \frac{1}{2} \|x\|_1^2 - \frac{1}{2} \|x_0\|_1^2$  we have  $d_1(x) = \frac{1}{2} \|x - x_0\|_1^2$  for  $x \in \mathcal{L}$ .

In an analogous fashion we define the center  $y_0$  of  $Q_2$  and assume that  $d_2$  vanishes at  $y_0$ . Therefore

$$d_2(y) \geq \frac{1}{2} \sigma_2 \|y - y_0\|_2^2.$$

Finally, let  $D_1$  and  $D_2$  satisfy  $D_1 \geq \max_x \{d_1(x) : x \in Q_1\}$  and  $D_2 \geq \max_y \{d_2(y) : y \in Q_2\}$ .

PROPOSITION 14 (Nesterov [13, Theorem 1]). *For  $\mu > 0$ , the function*

$$(5.7) \quad \varphi_\mu(x) := \max_y \{ \langle Ax, y \rangle - \mu d_2(y) : y \in Q_2 \}$$

*is a continuously differentiable uniform approximation of  $\varphi$ :*

$$(5.8) \quad \varphi_\mu(x) \leq \varphi(x) \leq \varphi_\mu(x) + \mu D_2 \quad \text{for all } x \in \mathbf{E}_1.$$

*Moreover, if we denote by  $y_\mu(x)$  the (unique) maximizer from (5.7), then the gradient of  $\varphi_\mu(x)$  is given by  $\nabla \varphi_\mu(x) = A^* y_\mu(x)$  and is Lipschitz continuous with constant*

$$(5.9) \quad \gamma_\mu = \frac{1}{\mu \sigma_2} \|A\|_{1,2}^2.$$

The smooth version of (P) therefore is

$$(P'_{\text{sm}}) \quad \min_x \{ \varphi_\mu(x) : x \in Q_1 \}.$$

The main result of [13] is the following theorem.

THEOREM 15 (Nesterov [13, Theorem 3]). *If we apply Algorithm 6 to problem (P'\_{sm}) with smoothing parameter*

$$(5.10) \quad \mu = \frac{2\|A\|_{1,2}}{N+1} \sqrt{\frac{D_1}{\sigma_1 \sigma_2 D_2}}$$

*and if  $x = \text{Smooth}(\varphi_\mu, \gamma_\mu, Q_1, x_0, N)$ , then<sup>3</sup>*

$$(5.11) \quad \varphi(x) - \varphi^* \leq \frac{4\|A\|_{1,2}}{N+1} \sqrt{\frac{D_1 D_2}{\sigma_1 \sigma_2}}.$$

---

<sup>3</sup>The original theorem states the result as a gap between  $\varphi(x)$  and  $\theta(y)$  for a certain  $y \in Q_2$ .

**ALGORITHM 6. (SMOOTH) EFFICIENT SMOOTH METHOD.**


---

```

1: Input:  $f, \gamma, Q_1, x_0, N$ ;
2: for  $k = 0$  to  $N$  do
3:   Compute  $\nabla f(x_k)$ ;
4:    $y_k = \arg \min\{\langle \nabla f(x_k), x - x_k \rangle + \frac{\gamma}{2} \|x - x_k\|_1^2 : x \in Q_1\}$ ;
5:    $z_k = \arg \min\{\sum_{i=0}^k \frac{i+1}{2} \langle \nabla f(x_i), x - x_i \rangle + \frac{\gamma}{\sigma_1} d_1(x) : x \in Q_1\}$ ;
6:    $x_{k+1} = \frac{2}{k+3} z_k + \frac{k+1}{k+3} y_k$ ;
7: end for
8: Output:  $y_N$ 

```

---

**5.3. The main result.** We will use the above theorem in the same way as Proposition 3 to devise a  $O(1/\delta)$ -algorithm for finding a  $\delta$ -approximate solution of  $(P)$ . Algorithms of this type, formulated for several specific choices of objective functions, were proposed already by Nesterov [14], [15]. These methods are similar in spirit to Algorithm 2, recursively updating an upper bound on  $\varphi^*$ . We give an improved version of this algorithm applicable to the problems considered in the cited papers. Our contribution lies mainly in improving the theoretical complexity by incorporating a bisection speedup. As in the previous section, it is possible to formulate a nonrestarting version of our algorithm by sacrificing the double logarithm in the theoretical complexity for a single one.

**5.3.1. Preliminaries.** Let us return to problem  $(P)$ , using the representation (5.2) for the objective function (hence  $Q_1 = \mathcal{L}$ ), and approach it with the tools described in the previous subsections. Let  $\mathbf{E}_1 := \mathbf{E}$  and assume that  $G: \mathbf{E}_1 \rightarrow \mathbf{E}_1^*$  defines an ellipsoidal rounding of  $\partial\varphi(0) = A^*Q_2$  such that (3.1) holds. Notice that the inequalities (3.5), (3.6), and (3.7) are implied by (3.1). To be able to obtain an algorithm guaranteeing a  $\delta$ -approximate solution in relative scale, it is crucial to choose  $\|x\|_1 \equiv \|x\|_G, x \in \mathbf{E}_1$ .

If we wish to apply Algorithm 6, we need to supply to it a *bounded* subset of  $\mathcal{L}$  containing the minimizer. Observe that as long as we are in the possession of an upper bound  $R$  on  $\varphi^*$ , (3.7) guarantees that all minimizers of  $(P)$  lie in the set

$$Q_1(R) \stackrel{\text{def}}{=} \mathcal{L} \cap \{x: \|x - x_0\|_G \leq R\}.$$

The point  $x_0$ —the projection of the origin onto  $\mathcal{L}$  in the  $G$ -norm—is the center of  $Q_1(R)$  as defined in (5.6) if we choose the prox-function for  $Q_1(R)$  to be

$$d_1(x) := \frac{1}{2} \|x - x_0\|_G^2.$$

In this case  $\sigma_1 = 1$  and  $D_1 = \max\{d_1(x): x \in Q_1(R)\} = \frac{1}{2} R^2$ . We leave the choice of  $d_2$  purposely open to allow for fine tuning for particular applications.

A direct consequence of Theorem 15 with the settings described above is the following analogue of Lemma 5.

LEMMA 16. *If  $\varphi^* \leq R$ ,  $\beta > 0$ , and we set*

$$N = \left\lceil \frac{2\sqrt{2}\|A\|_{1,2}}{\beta} \sqrt{\frac{D_2}{\sigma_2}} \right\rceil, \quad \mu = \frac{\sqrt{2}\|A\|_{1,2}R}{N+1} \sqrt{\frac{1}{\sigma_2 D_2}},$$

and  $\gamma_\mu$  as in (5.9), then

$$x = \text{Smooth}(\varphi_\mu, \gamma_\mu, Q_1(R), x_0, N)$$

satisfies

$$\varphi(x) - \beta R \leq \varphi^* \quad \text{and} \quad \varphi(x) \leq (1 + \beta)R.$$

The above lemma leads to a bisection algorithm (Algorithm 7) in the same way as we have seen it in the section on subgradient algorithms. The main result follows.

---

**ALGORITHM 7. (SMOOTHBIS) SMOOTH BISECTION SCHEME.**

- 1: **Input:**  $\varphi, x_0, \rho, \beta, \tau, \delta$ ;
  - 2:  $k = 0, x = x_0, L_0 = \|x_0\|_G, R_0 = \varphi(x_0)$ ;
  - 3:  $c = (1 + \tau)(1 + \beta), N = \lfloor \frac{2\sqrt{2}\|A\|_{1,2}}{\beta} \sqrt{\frac{D_2}{\sigma_2}} \rfloor$ ;
  - 4: **while**  $R_k/L_k > c$  **do**
  - 5:      $R = \sqrt{\frac{L_k R_k}{1 + \beta}}, \mu = \frac{\sqrt{2}\|A\|_{1,2} R}{N+1} \sqrt{\frac{1}{\sigma_2 D_2}}, \gamma_\mu = \frac{\|A\|_{1,2}^2}{\mu \sigma_2}$ ;
  - 6:      $x = \text{Smooth}(\varphi_\mu, \gamma_\mu, Q_1(R), x_0, N)$ ;
  - 7:     **if**  $\varphi(x) \leq (1 + \beta)R$  **then**
  - 8:         set  $L_{k+1}, R_{k+1}$  as in (4.8)
  - 9:     **else**
  - 10:         set  $L_{k+1}, R_{k+1}$  as in (4.9)
  - 11:     **end if**
  - 12:      $k = k + 1$ ;
  - 13: **end while**
  - 14:  $N = \lfloor 2\sqrt{2} \frac{R_k}{L_k} \|A\|_{1,2} (1 + \frac{1}{\delta}) \sqrt{\frac{D_2}{\sigma_2}} \rfloor, \mu = \frac{\sqrt{2}\|A\|_{1,2} R_k}{N+1} \sqrt{\frac{1}{\sigma_2 D_2}}, \gamma_\mu = \frac{\|A\|_{1,2}^2}{\mu \sigma_2}$ ;
  - 15:  $x = \text{Smooth}(\varphi_\mu, \gamma_\mu, Q_1(R_k), x_0, N)$ ;
  - 16: **Output:**  $x$
- 

THEOREM 17. *Algorithm 7 returns a  $\delta$ -approximate solution of (P) and takes at most*

$$(5.12) \quad 2\sqrt{2}\|A\|_{1,2} \sqrt{\frac{D_2}{\sigma_2}} \left[ \frac{1}{\beta} \left[ \log_2 \left( \frac{\ln \rho}{\ln(1 + \tau)} \right) \right] + (1 + \tau)(1 + \beta) \left( 1 + \frac{1}{\delta} \right) \right]$$

steps of the smooth optimization subroutine. If  $\beta$  and  $\tau$  are constants, this becomes

$$(5.13) \quad O\left( \|A\|_{1,2} \sqrt{\frac{D_2}{\sigma_2}} \left( \ln \ln \rho + \frac{1}{\delta} \right) \right).$$

A reasonable practical choice of the parameters  $\beta$  and  $\tau$  is

$$(5.14) \quad \beta = \sqrt{\delta}, \quad \tau = \frac{1}{2} \left( \sqrt{1 + \frac{4\beta}{\ln 2}} - 1 \right).$$

**5.4. A direct representation of the objective function.** We can get rid of the dependence on  $\|A\|_{1,2}$  in (5.13) by identifying  $\mathbf{E}_2$  with  $\mathbf{E}_1^*$  (and consequently  $\mathbf{E}_1$

with  $\mathbf{E}_2^*$ ). In this case we can simply choose  $A = I$  and consider the following structural model for the objective function:

$$\varphi(x) = \max_g \{\langle g, x \rangle : g \in Q_2\}.$$

Let us set  $\|g\|_2 = \|g\|_1^* = \|g\|_G^*$  and select the following prox-function for  $Q_2$  (with center at the origin):

$$d_2(g) = \frac{1}{2}(\|g\|_G^*)^2.$$

Clearly  $\sigma_2 = 1$  and  $D_2 \leq \frac{1}{2}\rho^2$ ; the second inequality follows from the ellipsoidal rounding inclusion (3.1). Also observe that since  $\|\cdot\|_2^* \equiv \|\cdot\|_1$ , we have

$$\|A\|_{1,2} = \max\{\|Ax\|_2^* : \|x\|_1 = 1\} = \max\{\|x\|_1 : \|x\|_1 = 1\} = 1.$$

Substituting the values of these parameters into (5.13) gives the complexity

$$O\left(\rho\left(\frac{1}{\delta} + \ln \ln \rho\right)\right).$$

*Remark 1.* Observe that, in principle, we do not lose generality by “excluding”  $A$  because we can simply set the “new”  $Q_2$  to be equal to the “old”  $A^*Q_2$ . However, this sacrifice in modeling flexibility means that  $Q_2$  always coincides with  $\partial\varphi(0)$ , which has to be of a simple structure for the algorithm to work efficiently. This is mainly due to the need to compute derivatives of  $\varphi_\mu$ , which amounts to solving (5.7)—a concave quadratic maximization problem over  $Q_2$ . If this problem can not be solved efficiently (say, in a closed form), the method will likely be impractical.

## 6. Scaling and complexity.

**6.1. Scaling.** It is natural to ask the following question, how do the “relative-scale” algorithms developed in this paper perform when we scale the objective function? Note that if we replace  $A$  by  $tA$  for some  $t > 0$  (effectively scaling  $\varphi$  by  $t > 0$ ), then (3.1) holds with  $G$  replaced by  $t^2G$ . Therefore, the inequalities in section 3.3 are valid, and so are all the results of this paper. Note that, in particular, the values of  $\rho$  and  $\|A\|_{1,2}$  remain unchanged. Looking at (4.4), (4.12), (4.14), (4.21), and (5.13), we see that the iteration complexities of the algorithms discussed in the paper are not affected by scaling.

**6.2. Complexity comparison.** Table 1 compares the iteration complexities of the algorithms discussed in this paper.

TABLE 1  
Summary of iteration complexities.

Method name	Algorithm #	Number of iterations
SubSearch	2	$O(\rho^2\delta^{-2} \ln \rho)$
SubBis	3	$O(\rho^2\delta^{-2} + \rho^2 \ln \ln \rho)$
SubSearchNR	4	$O(\rho^2\delta^{-2} \ln \rho)$
SubBisNR	5	$O(\rho^2\delta^{-2} + \rho^2 \ln \rho)$
SmoothBis	7	$O(\rho\delta^{-1} + \rho \ln \ln \rho)$

**7. Applications.** In this section we apply the fastest of the algorithms developed in this paper—the bisection algorithm based on smoothing **SmoothBis**—to several problems of the form (P).

**7.1. Minimizing the maximum of absolute values of linear functions.** In this subsection we consider problem (P) with the objective function from Example 3:

$$(7.1) \quad \min\{\varphi_\infty(x) : x \in \mathcal{L}\}.$$

Many seemingly unrelated problems can be reformulated in the above form. For example, by (7.1) one can model

- the truss topology design problem,
- the problem of the construction of a  $c$ -optimal statistical design, and
- the problem of finding a solution of an underdetermined linear system having the smallest  $\ell_1$  norm.

In all the examples above the feasible set  $\mathcal{L}$  is a hyperplane. We will now show how one can solve problem (7.1) using the results of section 5. A different approach for solving the problems above, simultaneously and in relative scale, was recently proposed by Richtárik [17], [19]. The iteration complexity is also  $O(\frac{1}{\delta})$ , but the approach uses very different techniques.

**7.1.1. Applying the algorithm.** We will work with the last of the three representations for the objective function from Example 3:

$$\varphi_\infty(x) = \max\{|\langle a_i, x \rangle| : i = 1, 2, \dots, m\} = \max_y \{\langle Ax, y \rangle : y \in Q_2\}$$

with  $Q_2$  being the unit simplex in  $\mathbf{R}^{2m}$  and  $A$  the  $2m \times n$  matrix with rows  $a_i, -a_i, i = 1, \dots, m$ . In addition, assume that the vectors  $a_i, i = 1, \dots, m$ , span  $\mathbf{E}_1^* = \mathbf{R}^n$ . It seems natural to choose  $\|y\|_2 := \sum_i |y^{(i)}|$  so that  $\|y\|_2 = 1$  for all  $y \in Q_2$ . If we let

$$d_2(y) := \ln 2m + \sum_{i=1}^{2m} y^{(i)} \ln y^{(i)}$$

and define  $0 \times \ln 0 := \lim_{\tau \downarrow 0} \tau \ln \tau = 0$ , then by the following lemma,  $d_2$  is a prox-function on  $Q_2$  with center  $y_0 := (\frac{1}{2m}, \dots, \frac{1}{2m})$ .

LEMMA 18. *The prox-function  $d_2$  is strongly convex on  $Q_2$ , with respect to  $\|\cdot\|_2$ , with convexity parameter  $\sigma_2 = 1$ .*

*Proof.* It suffices to show that  $d_2(y) \geq \frac{1}{2} \|y - y_0\|_2^2$ . This can be proved by elementary means using only the Cauchy–Schwarz inequality (see, e.g., Borwein and Lewis [4, Exercise 3.3.25(d)]) or, using differentiation and properties of convex functions, (Nesterov [13, Lemma 3]).  $\square$

It is easy to see that  $D_2 = \sup\{d_2(y) : y \in Q_2\} = \ln 2m$  (the supremum is attained at each of the boundary vertices). Finally, let us compute the norm of the linear operator  $A$ :

$$\|A\|_{1,2} = \max_{\|x\|_1=1} \|Ax\|_2^* = \max_{\|x\|_G=1} \|Ax\|_\infty = \max_{\|x\|_G=1} \varphi(x) \stackrel{(3.5)}{\leq} \rho.$$

In view of (5.9) we have  $\gamma_\mu \leq \frac{\rho^2}{\mu}$ . It is shown in Nesterov [13, Lemma 4] that

$$\varphi_\mu(x) = \mu \ln \left( \frac{1}{2m} \sum_{i=1}^m [e^{\langle a_i, x \rangle / \mu} + e^{\langle -a_i, x \rangle / \mu}] \right).$$

Since  $\partial\varphi(0) = \text{conv}\{\pm a_i : i = 1, 2, \dots, m\}$  is a centrally symmetric subset of  $\mathbf{R}^n$ , we may assume that a good rounding, with  $\rho = O(\sqrt{n})$ , is available to us. It can be computed efficiently in  $O(n^2 m \ln m)$  arithmetic operations. For details about algorithms we refer to [1], [10], [14], [22], [23].

**7.1.2. Complexity.** It follows from (5.13) that Algorithm 7 has the complexity

$$O\left(\sqrt{n \ln m} \left(\ln \ln n + \frac{1}{\delta}\right)\right).$$

This improves the result of Nesterov [14], where the author gives the bound

$$O\left(\frac{\sqrt{n \ln m}}{\delta} \ln n\right).$$

**7.2. Minimizing the sum of absolute values of linear functions.** Consider problem (P) with the following objective function:

$$\varphi_1(x) = \sum_{i=1}^m |\langle a_i, x \rangle|.$$

As usual, we assume that the vectors  $a_1, a_2, \dots, a_m$  span  $\mathbf{E}_1^*$ .

**7.2.1. Applying the algorithm.** Let  $\mathbf{E}_1 = \mathbf{E}_1^* = \mathbf{R}^n$  and  $\mathbf{E}_2 = \mathbf{E}_2^* = \mathbf{R}^m$ , and let us represent  $\varphi_1$  as

$$(7.2) \quad \varphi_1(x) = \max_y \{\langle Ax, y \rangle : y \in Q_2\},$$

where  $Q_2 = \{y \in \mathbf{R}^m : |y^{(i)}| \leq 1, i = 1, 2, \dots, m\}$  and  $A$  is the  $m \times n$  matrix with rows  $a_1, \dots, a_m$ . Usually we first find a rounding of  $\partial\varphi_1(0)$ , and using the rounding operator we define a norm on  $\mathbf{E}_1$ . Because of the simple structure of  $Q_2$ , we will instead start by defining  $\|y\|_2 := (\sum_i (y^{(i)})^2)^{1/2}$  and noting that this leads to a  $\sqrt{m}$ -rounding of  $Q_2$ :

$$(7.3) \quad \mathcal{B}(I, 1) \subseteq Q_2 \subseteq \mathcal{B}(I, \sqrt{m})$$

with  $I: \mathbf{R}^m \rightarrow \mathbf{R}^m$  denoting the identity operator. We will show now how this naturally leads to a rounding operator defined on  $\mathbf{E}_1$ , enjoying the same quality of rounding.

LEMMA 19 (Nesterov [15, Lemma 2]). *If the vectors  $a_1, \dots, a_m$  span  $\mathbf{R}^m$ , then  $\|x\|_1 := \|Ax\|_2^*$  defines a norm on  $\mathbf{R}^n$ . Moreover, if we let  $G := A^T A$  (a positive definite matrix), then  $\|\cdot\|_1 \equiv \|\cdot\|_G$  and  $\mathcal{B}(G, 1) \subseteq \partial\varphi(0) = A^T Q_2 \subseteq \mathcal{B}(G, \sqrt{m})$ .*

Let us define  $d_2(y) := \frac{1}{2} \|y\|_2^2$  so that the convexity parameter of this prox-function is  $\sigma_2 = 1$ . It follows from (7.3) that  $D_2 = \max\{d_2(y) : y \in Q_2\} \leq \frac{1}{2} m$ . Finally,

$$\|A\|_{1,2} = \max\{\|Ax\|_2^* : \|x\|_1 = 1\} = \max\{\|x\|_1 : \|x\|_1 = 1\} = 1.$$

**7.2.2. Complexity.** It follows from (5.13) that Algorithm 7 has the complexity

$$O\left(\sqrt{m}\left(\frac{1}{\delta} + \ln \ln m\right)\right).$$

This improves the following bound of Nesterov [15]:

$$O\left(\frac{\sqrt{m} \ln m}{\delta}\right).$$

**7.3. Minimizing the maximum of linear functions over a simplex.** The motivation for this problem is the computation of the value of a two-person zero-sum matrix game with nonnegative coefficients: Let  $\hat{A} \in \mathbf{R}^{m \times n}$  be a real matrix with nonnegative entries and rows  $a_1, \dots, a_m$ . Consider the following game. There are two players: a row player ( $R$ ) and a column player ( $C$ ). Player  $R$  chooses a probability distribution  $y$  over the rows of matrix  $\hat{A}$ , and  $C$  chooses a probability distribution  $x$  over the columns. After that,  $C$  pays  $y^T \hat{A} x$  dollars to  $R$ . Assume the players are *conservative*; that is,  $C$  wishes to minimize his worst-case loss and  $R$  wants to maximize his worst-case win. That is,  $C$  prefers to choose strategy

$$x^* \in \arg \min_{x \in \Delta_n} \max_{y \in \Delta_m} y^T \hat{A} x,$$

and, similarly,  $R$  wishes to choose strategy

$$y^* \in \arg \max_{y \in \Delta_m} \min_{x \in \Delta_n} y^T \hat{A} x.$$

The set  $\Delta_n$  (resp.,  $\Delta_m$ ) denotes the unit simplex in  $\mathbf{R}^n$  (resp.,  $\mathbf{R}^m$ ). A classical result by von Neumann [24] says that<sup>4</sup>

$$\varphi^* := \min_{x \in \Delta_n} \max_{y \in \Delta_m} y^T \hat{A} x = \max_{y \in \Delta_m} \min_{x \in \Delta_n} y^T \hat{A} x.$$

The value  $\varphi^*$  is called the *value of the game*. Note that if we let  $Q_1 := \Delta_n$  and

$$\varphi(x) = \max\{\langle a_i, x \rangle : i = 1, 2, \dots, m\},$$

then we can write  $\varphi^* = \min_x \{\varphi(x) : x \in Q_1\}$ .

**7.3.1. Applying the algorithm.** First observe that

$$\partial\varphi(0) = \text{conv}\{a_i : i = 1, 2, \dots, m\},$$

which fails to satisfy (2.1) due to the assumption on nonnegativity of the entries of  $\hat{A}$ . To remedy this situation, we will follow a trick suggested by Nesterov [14]. Notice that we are interested in  $\varphi$  as defined on  $\Delta_n$  only, which is a subset of the nonnegative orthant. Let us therefore define

$$\hat{\varphi}(x) \stackrel{\text{def}}{=} \max\{\langle a_i, |x| \rangle : i = 1, 2, \dots, m\},$$

---

<sup>4</sup>For a modern proof based on Fenchel duality, we refer to, for example, Exercise 4.2.16 in Borwein and Lewis [4].

where  $|x| = (|x_1|, \dots, |x_n|)$  and observe that  $\hat{\varphi}(x) = \varphi(x)$  for all  $x \in \mathbf{R}_+^n$ , and

$$\partial\hat{\varphi}(0) = \text{conv} \bigcup_{i=1}^m \{g: -a_i \leq g \leq a_i\}.$$

It is particularly interesting to note that  $\partial\hat{\varphi}(0)$  is a *sign-invariant set*, one that with every point  $g$  contains all points obtained by arbitrarily changing the signs of the coordinates of  $g$ . In fact,  $\partial\hat{\varphi}(0)$  is the smallest sign-invariant set containing  $\partial\varphi(0)$ . Nesterov shows that sign-invariant convex bodies admit a more efficient rounding algorithm than the more general centrally symmetric sets mainly due to the possibility of working only with *diagonal* positive definite matrices defining the rounding.

Instead of rounding  $\partial\varphi(0)$ , one can therefore find an ellipsoidal rounding of  $\partial\hat{\varphi}(0)$  (defined by a diagonal positive definite matrix  $G$ ) with  $\rho = O(\sqrt{n})$  and then deduce inequality (3.5), which holds for all  $x \in \mathbf{R}_+^n$  (Nesterov [14, Lemma 5]). The smoothing of  $\varphi$  (and hence of  $\hat{\varphi}$  on the domain of interest) can be performed in complete analogy with the situation in subsection 7.1. The choice of the representation of the objective function, the choice of the prox-function for  $Q_2$ , and the implied bounds are all identical (the only change is that the dimension drops from  $2m$  to  $m$ ).

**7.3.2. Complexity.** The iteration complexity of Algorithm 3 as applied to the problem of computing the value of a two-person matrix game with nonnegative coefficients is

$$O\left(\sqrt{n \ln m} \left(\frac{1}{\delta} + \ln \ln n\right)\right).$$

This improves the result of Nesterov [14, Algorithm 4.4], where the author gives the bound

$$O\left(\frac{\sqrt{n \ln m}}{\delta} \ln n\right).$$

**8. Computational experiments.** In this section we perform computational tests on problems of the structure described in section 7.1:

$$(8.1) \quad \min\{\varphi_\infty(x) \equiv \max\{|\langle a_i, x \rangle|, i = 1, \dots, m\} : \langle d, x \rangle = 1\}.$$

All experiments were done on a Windows XP desktop with Intel Core 2 Quad Q8300 CPU at 2.5 GHz with 3.46 GB of RAM. Algorithm 7 is run with constants  $\beta, \tau$  as given in (5.14). Rounding of the centrally symmetric set  $\partial\varphi_\infty(0) = \text{conv}\{\pm a_i, i = 1, \dots, m\}$  was in all cases done by Khachiyan's algorithm [9] with  $\rho = 1.1\sqrt{n}$ .

**8.1. Data: Truss topology design.** The data  $A = [a_1, \dots, a_m] \in \mathbf{R}^{n \times m}$  and  $d \in \mathbf{R}^n$  was generated using a formulation of the truss topology design (TTD) problem in the form (8.1). For details of the derivation we refer to [3, section 1.3.5] and [15], [19]. A brief description of the problem will suffice for our purposes: A two-dimensional (2D) rectangle of size  $a \times b$  is discretized into  $a \times b$  equidistant nodes. The  $a$  nodes "on the left" are attached to a wall, and a 2D force is applied at all the remaining nodes. There are a total of  $a(b-1)$  free nodes; the vector of forces is thus of dimension  $n = 2a(b-1)$ . In our formulation this vector is  $d$  (and always represents a single horizontal unit force applied at the right-middle node in the rightward direction). The TTD problem  $ttd(a, b)$  is the

TABLE 2  
*Initialization by ellipsoidal rounding ( $\varphi^* = 1$ ).*

Problem	$n$	$m$	$L_0 = \ x_0\ _G$	$R_0 = \varphi(x_0)$	$q_0 = R_0/L_0$	$\rho = 1.1\sqrt{n}$
<i>ttd</i> (3, 3)	12	28	0.4005	1.4188	3.5429	3.8105
<i>ttd</i> (5, 5)	40	200	0.4053	1.0266	2.5332	6.9570
<i>ttd</i> (7, 7)	84	748	0.3855	1.9134	4.9634	10.0817
<i>ttd</i> (9, 9)	144	2040	0.3717	2.3301	6.2690	13.2000
<i>ttd</i> (21, 5)	200	3332	0.5257	2.1363	4.0637	15.5563

problem of designing a structure of bars with endpoints in the nodes, with total weight of all the bars limited such that the total *compliance* of the truss is minimized. Compliance is a quantity proportional to the work performed by the system after the forces are applied until the nodes and bars are displaced to equilibrium. Dimension  $m$  represents the total number of potential bars. In all of the problems we allow any two nodes to be connected; overlapping bars are not allowed.

**8.2. Rounding.** Table 2 lists six *ttd*( $a, b$ ) test problems and quantities  $L_0, R_0, q_0$ , and  $\rho$  obtained after the initial rounding phase (i.e., computation of  $G$  satisfying (3.1)). For convenience, matrix  $A$  has in each case been scaled so that the optimal value of each problem is 1. This does not affect the algorithms (see section 6.1) and allows for straightforward comparison between methods that work in relative scale and absolute scale, due to the fact these two notions then coincide.

Note that for all problems  $\varphi(x_0)$  is already quite close to the optimal value. It is within 2.7% of optimum in the *ttd*(5, 5) case, and within a factor of 2.33 in the *ttd*(9, 9) case. We know from (4.2) that  $\varphi(x_0) \leq \rho$  must hold; the actual initial function values for our problems are much better than this bound. This suggests that the rounding stage does a very good job in preprocessing the problem for the optimization stage <sup>5</sup>.

**8.3. Nonrestarting versus restarting.** In this test we will contrast the behavior of the nonrestarting variant of the subgradient search scheme (**SubSearchNR**) against its restarting version (**SubSearch**). Let us equip both methods with the additional ability to quit the subgradient subroutine at step  $k$  in case a point  $x$  is found for which  $\varphi(x) < \varphi(\hat{x}_{k-1})/c$ . It is particularly interesting to see what happens for extremely small values of the decrease factor: we will choose  $c = 1 + 10^{-8}$ . This choice essentially means that the subgradient subroutine will be left immediately after a point is found which is better than the current best point (it also means that the theoretical complexity of both methods blows up). The nonrestarting method should have a clear advantage: it starts the next subgradient subroutine from the current best point, and hence it should not take too long before a new better point is found. In contrast, the restarting version starts the whole process again from  $x_0$ . Both methods will run their next call of the subgradient subroutine with smaller stepsizes. Method **SubSearchNR** with extremely small  $c$  can in view of (4.13) be interpreted as a subgradient method which adjusts its stepsize as soon as it gets new information about the distance of the current best point to the set of minimizers, which happens every time a new best point is found.

The results of this comparison, for the *ttd*(9, 9) problem, are given in Table 3. For both methods we list the number of lower-level subgradient iterations  $N$  it takes to

<sup>5</sup>For an  $O(1/\delta)$  algorithm in which the rounding and optimization stages coincide, see [19].

TABLE 3  
*Nonrestarting versus restarting:  $ttd(9, 9)$ .*

$\hat{\delta}$	SubSearchNR		SubSearch	
	$\delta = 0.1$	$\delta = 0.01$	$\delta = 0.1$	$\delta = 0.1$
10%	117 (3)	631 (6)	2,821 (69)	107,722 (433)
9%	134 (3)	703 (9)	3,666 (98)	122,493 (495)
8%	181 (31)	827 (15)	5,473 (116)	150,703 (592)
7%	207 (4)	1084 (12)	6,830 (167)	201,754 (862)
6%	223 (2)	1387 (9)	7,591 (213)	292,715 (1,125)
5%	420 (33)	1744 (5)	8,572 (278)	374,557 (1,468)

achieve a certain relative accuracy level  $\hat{\delta}$  (this is not the target accuracy  $\delta$  that enters the method as an input). The number in the parentheses is the number of subgradient steps in the last call of the subgradient subroutine, the one in which the  $\hat{\delta}$ -approximate point was found. The difference between the two methods is clear. For any given accuracy  $\hat{\delta}$ , the total number of subgradient steps of the nonrestarting method (**SubSearchNR**) is approximately equal to the number of subgradient steps of the last call of the subgradient subroutine in the restarting version (**SubSearch**). Notice that the choice of  $\delta$  has a huge effect on the performance of both methods, as it directly affects the stepsize of the subgradient subroutines. Smaller  $\delta$  leads to bigger  $N \sim (1 + \frac{1}{\delta^2})$ , which in turn leads to smaller stepsize  $\kappa = \frac{R}{\sqrt{N+1}}$ . A similar but milder effect occurs for the smooth methods as well: small  $\delta$  increases  $N$ , which decreases the smoothing parameter  $\mu$ , which increases the Lipschitz constant  $L_\mu$  of  $\nabla\varphi_\mu$ , which in turn leads to smaller steps via (8.3) ( $t = L_\mu/2$ ).

**8.4. Relativity speedup.** In this test we compare the fastest of our methods, **SmoothBis** (Algorithm 7), to **Smooth** (Algorithm 6) applied to the smoothed version of each problem *directly*, in the spirit of Theorem 15. We set  $\delta = 0.01$  for **SmoothBis** and  $\epsilon = 0.01$  for **Smooth** (these settings have the same meaning as  $\varphi^* = 1$  in all test problems). For each of the test problems and both algorithms we list the number of iterations  $N$  of the smooth subroutine (each comprising two subproblems of the form (8.2)), time  $t$  in seconds, and the accuracy  $\hat{\delta}$ ,  $\hat{\epsilon}$  at termination. Both algorithms are run for the full number of iterations, as prescribed by theory. The results are given in Table 4, and the bar structure of the resulting optimal trusses in Figure 2. Note that the method working in relative scale (**SmoothBis**) is faster on all test problems except  $ttd(3, 3)$ , both in terms of speed and iteration count. That is, we do not pay for obtaining a result in relative scale; quite to the contrary, we benefit from it.

TABLE 4  
*Relativity speedup.*

Problem	Smooth (Alg 6), $\epsilon = 0.01$			SmoothBis (Alg 7), $\delta = 0.01$		
	$N$	$t$	$\hat{\epsilon}$	$N$ (saving)	$t$ (saving)	$\hat{\delta}$
$ttd(3, 3)$	2,594	0.2''	$<10^{-20}$	2,990	0.3''	$<10^{-20}$
$ttd(5, 5)$	7,863	1.1''	$5.6 \times 10^{-4}$	6,030 (23.3%)	0.9'' (21.0%)	$4.9 \times 10^{-4}$
$ttd(7, 7)$	15,091	7.3''	$3.6 \times 10^{-4}$	9,344 (38.2%)	4.2'' (42.5%)	$3.2 \times 10^{-4}$
$ttd(9, 9)$	22,245	54.4''	$6.3 \times 10^{-4}$	13,053 (41.3%)	32.2'' (40.8%)	$5.7 \times 10^{-4}$
$ttd(21, 5)$	27,891	140.5''	$4.8 \times 10^{-4}$	15,961 (42.8%)	77.7'' (44.7%)	$4.3 \times 10^{-4}$

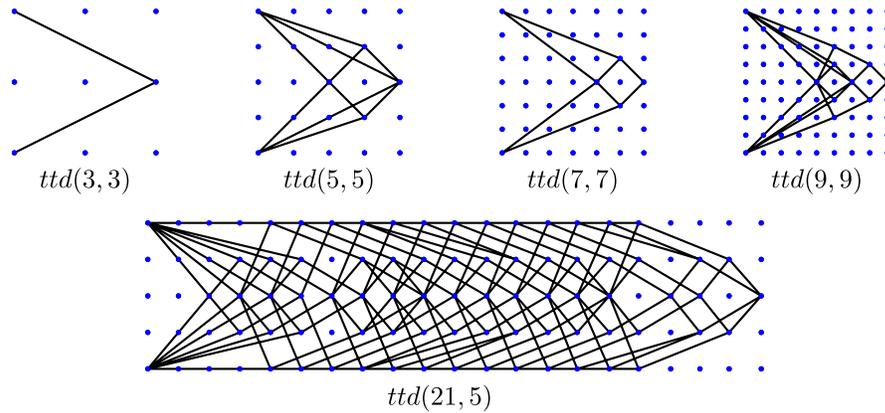


FIG. 2. Bar structure of optimal trusses ( $\delta = 0.001$ ).

**8.5. Bisection speedup.** Let us now compare the smooth bisection scheme, **SubBis** (Algorithm 7), with Algorithm 3.9 of Nesterov [14] (let us call it **SmoothSearch**)—a smooth analogue of Algorithm 2 in which the role of the subgradient subroutine **Subgrad** is replaced by **Smooth**. The iteration complexity of **SmoothSearch** is

$$\sqrt{8e\rho}(1 + \ln \rho)\sqrt{\ln(2m)}\left(1 + \frac{1}{\delta}\right)$$

with each step comprising of two operations of type (8.2). Table 5 compares the methods on a single problem,  $ttd(9, 9)$ , for several target relative accuracies  $\delta$ . Results similar to these were observed also for the other test problems, and we therefore do not list them. In particular, for each  $\delta$  and both methods, we list the number of iterations  $N$ , running time  $t$  in seconds, and the accuracy level  $\hat{\delta}$  at termination. Both methods are run for the full number of iterations as prescribed by the iteration complexity analysis. For the faster method (**SmoothBis**) we also list the percentage savings in iteration count and time as compared to the slower method. Notice that for both algorithms, the number of iterations increases linearly with decreasing  $\delta$ , as predicted by the theory. In all cases the termination accuracy is higher than the target accuracy by a bit more than an order of magnitude. Finally, observe that the advantage of **SmoothBis** grows, both in speed and number of iterations, with increasing accuracy demand.

TABLE 5  
Bisection speedup:  $ttd(9, 9)$ .

$\delta$	SmoothSearch			SmoothBis		
	$N$	$t$	real $\hat{\delta}$	$N$ (saving)	$t$ (saving)	real $\hat{\delta}$
0.05	6,145	18.1''	$2.2 \times 10^{-3}$	3,289 (46.5%)	12.3'' (32.0%)	$2.4 \times 10^{-3}$
0.01	29,555	68.9''	$4.8 \times 10^{-4}$	13,053 (55.8%)	33.7'' (51.1%)	$5.7 \times 10^{-4}$
0.005	58,818	125.0''	$2.4 \times 10^{-4}$	24,694 (58.0%)	57.3'' (54.2%)	$3.0 \times 10^{-4}$
0.001	292,919	575.3''	$4.8 \times 10^{-5}$	116,153 (60.4%)	225.9'' (60.7%)	$6.0 \times 10^{-5}$
0.0005	585,546	1078.1''	$2.4 \times 10^{-5}$	229,065 (60.9%)	440.7'' (59.1%)	$3.0 \times 10^{-5}$

**8.6. Algorithm 6.** Steps 4–5 of Algorithm 6 are of the form

$$(8.2) \quad \min\{\langle s, x \rangle + t\|x - \bar{x}\|_G^2 : \langle d, x \rangle = 1, \|x - x_0\|_G \leq R\},$$

where  $\langle d, \bar{x} \rangle = 1$  and  $t > 0$ . The solution is given by

$$(8.3) \quad x = x_0 - \frac{1}{2(t + \alpha)} G^{-1}(s' + \lambda d), \quad \text{where}$$

$$s' = s + 2tG(x_0 - \bar{x}), \quad \lambda = \langle s', x_0 \rangle, \quad v = \|s' + \lambda d\|_G^*, \quad \alpha = \begin{cases} \frac{v}{2R} - t, & v \geq 2Rt, \\ 0, & \text{otherwise.} \end{cases}$$

**Acknowledgment.** The author is very grateful to Mike Todd for numerous enlightening discussions and encouragement to publish these results.

#### REFERENCE

- [1] D. AHİPAŞAOĞLU, P. SUN, AND M. J. TODD, *Linear convergence of a modified Frank–Wolfe algorithm for computing minimum-volume enclosing ellipsoids*, *Optim. Methods Softw.*, 23 (2008), pp. 5–19.
- [2] A. BELLONI AND R. M. FREUND, *On the symmetry function of a convex set*, *Math. Program.*, 111 (2007), pp. 57–93.
- [3] A. BEN-TAL AND A. NEMIROVSKI, *Lectures on Modern Convex Optimization: Analysis, Algorithms, and Engineering Applications*, SIAM, Philadelphia, 2001.
- [4] J. M. BORWEIN AND A. S. LEWIS, *Convex Analysis and Nonlinear Optimization: Theory and Examples*, CMS Books Math., Springer-Verlag, New York, 2000.
- [5] F. A. CHUDAK AND V. ELEUTÉRIO, *Improved approximation schemes for linear programming relaxations of combinatorial optimization problems*, in *Integer Programming and Combinatorial Optimization*, Lecture Notes in Comput. Sci. 3509, M. Jünger and V. Kaibel, eds., Springer-Verlag, Berlin, 2005, pp. 191–219.
- [6] J.-L. GOFFIN, *On convergence rates of subgradient optimization methods*, *Math. Program.*, 13 (1977), pp. 329–347.
- [7] J.-B. HIRIART-URRUTY AND C. LEMARÉCHAL, *Convex Analysis and Minimization Algorithms*, Springer-Verlag, Berlin, 1993.
- [8] F. JOHN, *Extremum problems with inequalities as subsidiary conditions*, in *Studies and Essays, Presented to R. Courant on his 60th Birthday January 8, 1948*, Wiley Interscience, New York, 1948, pp. 187–204.
- [9] L. G. KHACHIYAN, *Rounding of polytopes in the real number model of computation*, *Math. Oper. Res.*, 21 (1996), pp. 307–320.
- [10] P. KUMAR AND E. A. YILDIRIM, *Minimum volume enclosing ellipsoids and core sets*, *J. Optim. Theory Appl.*, 126 (2005), pp. 1–21.
- [11] YU. NESTEROV, *A method for unconstrained convex minimization problem with the rate of convergence  $O(\frac{1}{k^2})$* , *Soviet. Math. Doct.*, 269 (1983), pp. 543–547.
- [12] YU. NESTEROV, *Introductory Lectures on Convex Optimization. A Basic Course*, Applied Optimization 87, Kluwer, Boston, 2004.
- [13] YU. NESTEROV, *Smooth minimization of non-smooth functions*, *Math. Program.*, 103 (2005), pp. 127–152.
- [14] YU. NESTEROV, *Rounding of convex sets and efficient gradient methods for linear programming problems*, *Optim. Methods Softw.*, 23 (2008), pp. 109–128.
- [15] YU. NESTEROV, *Unconstrained convex minimization in relative scale*, *Math. Oper. Res.*, 34 (2009), pp. 180–193.
- [16] YU. NESTEROV, *Barrier subgradient method*, *Math. Program.*, 127 (2011), pp. 31–56.
- [17] P. RICHTÁRIK, *Some Algorithms for Large-Scale Linear and Convex Minimization in Relative Scale*, Ph.D. thesis, Cornell University, School of Operations Research and Information Engineering, Ithaca, NY, 2007.
- [18] P. RICHTÁRIK, *Approximate level method for nonsmooth convex minimization*, *J. Optim. Theory Appl.*, 152 (2012), DOI [10.1007/s10957-011-9908-1](https://doi.org/10.1007/s10957-011-9908-1).
- [19] P. RICHTÁRIK, *Simultaneously Solving Seven Optimization Problems in Relative Scale*, manuscript, 2008.

- [20] R. T. ROCKAFELLAR, *Convex Analysis*, Princeton University Press, Princeton, NJ, 1997.
- [21] N. Z. SHOR, *Minimization Methods for Nondifferentiable Functions*, Springer-Verlag, Berlin, 1985.
- [22] P. SUN AND R. M. FREUND, *Computation of minimum-volume covering ellipsoids*, *Oper. Res.*, 52 (2004), pp. 690–706.
- [23] M. J. TODD AND E. A. YILDIRIM, *On Khachiyan's algorithm for the computation of minimum volume enclosing ellipsoids*, *Discrete Appl. Math.*, 155 (2007), pp. 1731–1744.
- [24] J. VON NEUMANN AND O. MORGENSTERN, *The Theory of Games and Economic Behavior*, Princeton University Press, Princeton, NJ, 1948.