

# Simultaneously solving seven optimization problems in relative scale

Peter Richtárik\*

December 2008

## Abstract

In this paper we develop and analyze an efficient algorithm which solves seven related optimization problems simultaneously, in relative scale. Each iteration of our method is very cheap, with main work spent on matrix-vector multiplication. We prove that if a certain sequence generated by the algorithm remains bounded, then the method must terminate in  $O(1/\delta)$  iterations, producing  $\delta$ -approximate solutions to all seven problems, where  $\delta$  is the prescribed relative error.

The seven problems are 1) a specific sublinear minimization program with a single homogeneous linear constraint, 2-3) the problem of finding the intersection of a ray and a centrally-symmetric polytope represented as a convex hull of a collection of points, 4) centrally-symmetric linear programming, 5) the problem of finding the least  $\ell_1$ -norm solution of an underdetermined linear system (basis pursuit), 6) a certain smooth convex minimization problem on the unit simplex and, finally, 7) a semidefinite program with rank-one objective and constraint matrices.

We finish the discussion by describing applications to truss topology design and design of statistical experiments.

**Keywords:** nonsmooth convex optimization, variable-metric subgradient method, linear programming (LP), centrally symmetric linear programming (CSLP), ellipsoid-squeezing method, sparse solution of underdetermined linear systems, iteratively reweighted least squares (IRLS), iteratively reweighted Euclidean projection (IREP), basis pursuit,  $\ell_1$ -norm minimization,  $c$ -optimal design, Elfving set, multiplicative-update method, Khachiyan's ellipsoidal rounding algorithm, Frank-Wolfe algorithm.

**Mathematics Subject Classification (2000):** 65K05, 62K05, 90C05, 90C06, 90C25.

---

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

The results of this paper were obtained in the years 2006 and 2007 while the author was a research assistant at Cornell University, working under the guidance of Mike Todd. This research was partially supported by NSF through grants DMS-0209457 and DMS-0513337 and by ONR through grant N00014-02-0057.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Seven problems: overview . . . . .	4
1.2	Contents . . . . .	6
1.3	Notation and basic concepts . . . . .	6
<b>2</b>	<b>Seven optimization problems</b>	<b>8</b>
2.1	The first five problems . . . . .	8
2.2	The main problem . . . . .	10
2.3	The seventh problem . . . . .	12
2.4	First algorithmic idea . . . . .	13
<b>3</b>	<b>Convexity and smoothness</b>	<b>13</b>
3.1	Convexity of the domain . . . . .	14
3.2	Convexity of the objective function . . . . .	15
3.3	Smoothness . . . . .	16
3.4	An example . . . . .	16
<b>4</b>	<b>Optimality conditions</b>	<b>17</b>
<b>5</b>	<b>Rank-one update</b>	<b>21</b>
5.1	A multiplicative weight update algorithm . . . . .	21
5.2	Ingredients of a rank-one update algorithm . . . . .	22
<b>6</b>	<b>Line search</b>	<b>27</b>
6.1	General line-search formula . . . . .	28
6.2	Specialized line-search formula . . . . .	29
<b>7</b>	<b>An algorithm with “increase” steps only</b>	<b>31</b>
7.1	A step-size heuristic . . . . .	31
7.2	Sufficient decrease . . . . .	32
7.3	Even better decrease . . . . .	33
7.4	A crucial assumption . . . . .	33
7.5	Quick and dirty analysis . . . . .	33
7.6	Refined analysis . . . . .	35
<b>8</b>	<b>An algorithm with “increase”, “decrease” and “drop” steps</b>	<b>36</b>
<b>9</b>	<b>Bounding the unknown constant</b>	<b>39</b>
9.1	Bounding the weights away from zero: theoretical implications . . . . .	40
9.2	Bounding the weights away from zero: algorithmic implications . . . . .	41
9.3	The average of the gammas . . . . .	42

9.4	An alternative proof of the bound on $\gamma_j$ . . . . .	43
<b>10</b>	<b>Interpretation</b>	<b>43</b>
10.1	(P3): The Frank-Wolfe algorithm on the unit simplex . . . . .	43
10.2	(P2): An ellipsoid-squeezing method for centrally-symmetric LP . . . . .	44
10.3	(D2): An Iteratively Reweighted Euclidean Projection Algorithm . . . . .	45
10.3.1	Two remarks . . . . .	46
10.3.2	The first iterate . . . . .	46
<b>11</b>	<b>Applications</b>	<b>47</b>
11.1	Truss topology design . . . . .	47
11.1.1	The problem . . . . .	47
11.1.2	Correspondence with the setting of problem (P3) . . . . .	47
11.1.3	Three examples . . . . .	48
11.2	Optimal design of statistical experiments . . . . .	49

## 1 Introduction

In this paper we propose and analyze the *iteration complexity* of two variants of an algorithm for *simultaneously* and *efficiently* solving several *large-scale* structured optimization problems in *relative scale*.

For the complexity results we need to assume that certain values generated by the algorithm stay bounded (Assumption 33). Under this condition we prove that the methods need at most  $O(1/\delta)$  iterations to produce a solution within relative error  $\delta$ . Each iteration of our method is very cheap — we only need to perform matrix-vector multiplications — and hence the proposed approach is suitable for solving very large instances.

By relative scale we mean that the guaranteed bound on the error is not additive, but multiplicative, relative to the optimal objective value. For example, for a minimization problem with objective function  $\varphi$  and optimal value  $\varphi^* > 0$ , we accept feasible points  $x$  satisfying the inequality

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

Algorithms for continuous optimization problems which guarantee results in relative scale, and are hence insensitive to the scaling of the objective function, are not common in the literature. For some recent work on this topic we refer the reader to, for example, [21, 22, 24, 25].

Primal-dual methods — algorithms that solve both a primal and a dual problem at the same time — are common in the optimization literature. In this work we can solve more problems because of the simple underlying structure allowing for several interesting reformulations.

## 1.1 Seven problems: overview

Our algorithm simultaneously solves three pairs of primal and dual problems:

$$(P1) - (D1), \quad (P2) - (D2), \quad \text{and} \quad (P3) - (D3),$$

and, additionally, an “outer” variant ( $D'1$ ) of ( $D1$ ). Their formal definitions are given and their interrelatedness is outlined in Section 2. It is interesting that our basic method, as applied to any of these problems, has a *natural interpretation* in that particular setting. We will now proceed with a brief informal overview of the various problems and comment on the characteristics of our algorithm in each case.

1. *Problem (P1)*: We start the discussion with a minimization program with a single linear equality constraint and objective function which is the maximum of a finite number of absolute values of homogeneous linear functions. Note that such a function is necessarily sublinear (convex and homogeneous of degree one), sign-symmetric and nonnegative.

*Interpretation:* One variant of our method can be interpreted as a variable-metric sub-gradient method [27] with explicit update rules for the metric and step lengths. The iteration-complexity is improved by an order of magnitude from the black-box optimal value of  $O(1/\delta^2)$  [17] and matches the result obtained by a general smoothing technique developed by Nesterov [20]. A preprocessing step in Nesterov’s approach [22] for solving problem ( $P1$ ) in relative scale within  $O(1/\delta)$  iterations involves the computation of an approximate ellipsoidal rounding of the underlying polytope, explicit smoothing of the nonsmooth objective function, and applying an optimal smooth method [18–20] to the new smoothed problem. In contrast, we utilize the very special structure of this particular problem by modifying the rounding procedure so that it becomes, at the same time, the optimization procedure. This is a way how to circumvent smoothing and still obtain an improved complexity analysis.

2. *Problem (D1)*: In the second problem, which is equivalent to the dual of ( $P1$ ), we are trying to find the intersection of a ray and a centrally-symmetric polytope arising as a convex hull of line segments with centers in the origin. This polytope is the subdifferential at the origin of the objective function of the previous problem. The feasible region is the line segment joining the origin and the intersection point.

*Interpretation:* Our method generates a sequence of ellipsoid with centers in the origin, inscribed in the polytope, and cutting the ray at points ever closer to the searched-after intersection point. The ellipsoids are updated in a rank-one fashion in a way similar to Khachiyan’s algorithm for computing a pair of Löwner-John ellipsoids [15, 21]. The goal at every iteration is no longer to maximize the volume of the next ellipsoid, but to cut a “larger” portion of the ray by the next ellipsoid. We show in Section 6 that a closed-form formula can be obtained for the optimal update.

3. *Problem (D'1)*: This is the same problem as (D1), except feasible points are “outside” the polytope and not “inside”. In other words, we are searching the portion of the ray which does not lie in the polytope.

*Interpretation:* It turns out that our algorithm also generates a sequence of points sitting on the portion of the ray which is outside of the polytope, and converging to the intersection point.

4. *Problem (P2)*: This is the problem of maximizing a linear function over a centrally symmetric polytope defined by pairs of cutting planes — the polar set of the polytope from (D1). Hence this is a *centrally-symmetric linear programming problem* (CSLP).

*Interpretation:* Our algorithm produces a sequence of ellipsoids centered in the origin containing the feasible set. These are the polar ellipsoids to those described (D1). New ellipsoid is obtained by *squeezing* the previous ellipsoid along the direction of the objective vector, with the goal to make it as thin as possible in this direction, subject to the constraint that only rank-one update squeezes are admissible. We will refer to our algorithm in the setting of CSLP as the *ellipsoid-squeezing method*. More detail can be found in Section 10.2.

5. *Problem (D2)*: Here we are asked to compute the  $\ell_1$ -projection of the origin onto an affine space given by an underdetermined set of linear equations. This is also known as the *Basis Pursuit* model for the NP-hard [9, 16] problem of finding sparse solutions of linear systems [5–7]. Problem (D2) is dual to (P2).

*Interpretation:* Our algorithm produces a sequence of Euclidean projections, with diagonal metrics, converging to the desired point. In spirit, this approach is similar to existing algorithms in the overdetermined case, known as iteratively reweighted least squares (IRLS). For more information on IRLS we refer the reader to the recent paper [8] and the references therein. In analogy, it seems natural to refer to our method in this setting as *iteratively reweighted Euclidean projection* (IREP). We discuss this problem in more detail in Section 10.3.

6. *Problem (P3)*: This is center of our focus in this paper. All analysis and presentation is done for this particular problem and in Section 3 we show that an approximate solution to it leads to approximate solutions of all other problems. In particular, (P3) is a convex minimization problem on the unit simplex with a specific structured objective function which is smooth on the interior of the simplex but can have infinite values on portions of the boundary.

In statistics, the problem precisely corresponds to computing the *c-optimal design* over a finite regression space. Various algorithms have been proposed in the statistics literature for solving it [10, 23, 28]. Our method falls into the category of *multiplicative update algorithms*. To the best of our knowledge, no iteration complexity results are known in the statistics literature for this problem.

*Interpretation:* Our algorithm can be interpreted as the Frank-Wolfe method (with or without “away” steps) [11] with exact line-search formulae. We will discuss this in a bit more detail in Section 10.1.

7. *Problem (D3):* The last problem is a semidefinite program with a rank-one objective matrix and several rank-one inequality constraint matrices.

*Interpretation:* Our algorithm essentially solves this problem by utilizing the fact that there must be an optimal rank-one solution (which also follows from Theorem 21), which transforms the problem into (P2). From the perspective of design of experiments, the relationship between (P3) and (D3) is known as mutual boundedness theorem for scalar optimality [23].

## 1.2 Contents

The paper is organized as follows. In Section 2 we formally describe all seven problems and explore the relationships among them. Having done that, we then establish convexity and smoothness of the objective function of the main problem, (P3), in Section 3. Section 4 contains a result stating that a certain single condition implies approximate optimality in all seven problems. Rank-one updates and their properties needed for further analysis are discussed in Section 5. In Section 6 we calculate an explicit exact line-search formula for (P3). The next two sections contain the core of the paper: the algorithms and their iteration-complexity analysis. In Section 9 we comment on the technical “boundedness” assumption. Next, in Section 10 we offer a more detailed discussion of our algorithm as applied to three of the problems, expanding on the outline contained in the introduction. Finally, in the last section we describe two of the many possible applications of our methods: to truss topology design and to computing *c*-optimal designs of statistical experiments.

## 1.3 Notation and basic concepts

Throughout the paper,  $\mathbf{E}$  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 primary objects of this paper are nonzero vectors  $d, a_1, \dots, a_m \in \mathbf{E}^*$  and the centrally symmetric convex set

$$Q \stackrel{\text{def}}{=} \text{conv}\{\pm a_i : i = 1, 2, \dots, m\}. \tag{1}$$

By  $\langle g, x \rangle$  we denote the action of the functional  $g \in \mathbf{E}^*$  on  $x \in \mathbf{E}$ . If  $G : \mathbf{E} \rightarrow \mathbf{E}^*$  is a linear operator,  $\text{null}(G)$ ,  $\text{span}(G)$ ,  $\text{range}(G)$  denote its nullspace, span and range, respectively. If we treat  $G$  as a matrix (after fixing a pair of bases),  $\text{rank}(G)$  denotes the rank of  $G$ ,  $\text{diag}(G)$  is a vector containing the diagonal of  $G$  and  $\text{Diag}(w)$  is a diagonal matrix with vector  $w$  on its diagonal. Coordinates of a vector are denoted by superscripts in brackets; for example,  $w = (w^{(1)}, \dots, w^{(m)})$ . Subscript is used to enumerate scalar (mostly Greek letters) or vector variables. For the unit simplex in  $\mathbf{R}^m$  we use the notation

$$\Delta_m \stackrel{\text{def}}{=} \left\{ w \in \mathbf{R}^m : \sum_i w^{(i)} = 1, w^{(i)} \geq 0 \right\}.$$

For  $v \in \mathbf{R}^m$  we write  $\|v\|_1 = \sum |v^{(i)}|$  and  $\|v\|_\infty = \max |v^{(i)}|$ .

The *support function* of a nonempty set  $\mathcal{X} \subset \mathbf{E}$  is the function  $\xi_{\mathcal{X}} : \mathbf{E}^* \rightarrow \bar{\mathcal{R}}$  defined by

$$\xi_{\mathcal{X}}(g) \stackrel{\text{def}}{=} \sup\{\langle g, x \rangle : x \in \mathcal{X}\}.$$

For example,  $\xi_Q(x) = \max\{|\langle a_i, x \rangle| : i = 1, 2, \dots, m\}$ . The *polar* of a convex set  $\mathcal{X} \subset \mathbf{E}$  is the set  $\mathcal{X}^\circ \subset \mathbf{E}^*$  defined by

$$\mathcal{X}^\circ \stackrel{\text{def}}{=} \{g \in \mathbf{E}^* : \langle g, x \rangle \leq 1 \text{ for all } x \in \mathcal{X}\}.$$

For example,

$$Q^\circ = \{x \in \mathbf{E} : |\langle a_i, x \rangle| \leq 1, i = 1, 2, \dots, m\}. \quad (2)$$

Analogous definitions work for objects defined in  $\mathbf{E}^*$ .

For a positive semidefinite self-adjoint linear operator  $G : \mathbf{E} \rightarrow \mathbf{E}^*$  define

$$\|x\|_G \stackrel{\text{def}}{=} \langle Gx, x \rangle^{1/2}. \quad (3)$$

It can easily be verified that

$$\|x\|_G = 0 \quad \Leftrightarrow \quad x \in \text{null}(G). \quad (4)$$

Also let

$$\|g\|_G^* \stackrel{\text{def}}{=} \begin{cases} \langle g, x \rangle^{1/2} & \text{if } g \in \text{range}(G) \text{ with } Gx = g, \\ +\infty & \text{otherwise.} \end{cases} \quad (5)$$

Notice that  $\langle g, x_1 \rangle = \langle g, x_2 \rangle$  whenever  $Gx_1 = g$  and  $Gx_2 = g$  because  $G$  is self-adjoint and hence  $\langle g, x_1 \rangle = \langle Gx_2, x_1 \rangle = \langle Gx_1, x_2 \rangle = \langle g, x_2 \rangle$ , all of which are nonnegative since  $G \succeq 0$  and, for example,  $\langle g, x_1 \rangle = \langle Gx_1, x_1 \rangle \geq 0$ . Hence (5) gives a valid definition.

Let us establish special notation for the sublevel sets of  $\|\cdot\|_G$  and  $\|\cdot\|_G^*$ :

$$\mathcal{B}^\circ(G) \stackrel{\text{def}}{=} \{x \in \mathbf{E} : \|x\|_G \leq 1\}, \quad \text{and} \quad (6)$$

$$\mathcal{B}(G) \stackrel{\text{def}}{=} \{g \in \mathbf{E}^* : \|g\|_G^* \leq 1\}. \quad (7)$$

Note that  $\mathcal{B}^\circ(G)$  is an *ellipsoidal cylinder* in  $\mathbf{E}$  and  $\mathcal{B}(G)$  is an *ellipsoid* in  $\text{range}(G)$ .

## 2 Seven optimization problems

### 2.1 The first five problems

For  $x \in \mathbf{E}$  let

$$\varphi(x) \stackrel{\text{def}}{=} \xi_Q(x) = \max_i |\langle a_i, x \rangle|, \quad (8)$$

and consider the problem

$$\boxed{\varphi^* \stackrel{\text{def}}{=} \min_x \{\varphi(x) : \langle d, x \rangle = 1\}.} \quad (P1)$$

The objective function is a nonnegative sublinear (convex and positively homogeneous) function with subdifferential at the origin equal to  $Q$ . Note that  $0 \in Q$ . We will however further assume that

$$0 \in \text{int } Q. \quad (9)$$

This implies that  $\varphi$  vanishes only at the origin, which is then the *unique* global minimizer of  $\varphi$ , whence  $\varphi^* > 0$ . Assumption (9) is equivalent to

$$\text{range}(A) = \mathbf{E}^*, \quad (10)$$

where  $A = [a_1, \dots, a_m]: \mathbf{R}^m \rightarrow \mathbf{E}^*$  is the linear operator mapping the  $i$ -th unit vector of  $\mathbf{R}^m$  to  $a_i$ . By  $A^*$  we denote the adjoint of  $A$ . This is the operator  $A^*: \mathbf{E} \rightarrow (\mathbf{R}^m)^*$  defined by  $\langle Av, x \rangle = \langle A^*x, v \rangle$  for all  $x \in \mathbf{E}$ ,  $v \in \mathbf{R}^m$ , so that  $A^*x = [\langle a_1, x \rangle, \dots, \langle a_m, x \rangle]^T$  and hence  $\varphi(x) = \|A^*x\|_\infty$ .

The (Lagrangian) dual of problem (P1) can be shown to be equivalent to

$$\boxed{\varphi^* = \max_\tau \{\tau : \tau d \in Q\},} \quad (D1)$$

and thus

$$\varphi^* d \in \text{bdry } Q. \quad (11)$$

As an exercise, let us check weak duality. Assume we have  $x$  with  $\langle d, x \rangle = 1$  and  $\tau$  with  $\tau d \in Q$ . Then  $\tau d$  is a weighted average of points from  $\{\pm a_i, i = 1, 2, \dots, m\}$  and hence  $\tau = \langle \tau d, x \rangle$  is equal to a weighted average of inner products from  $\{\langle \pm a_i, x \rangle, i = 1, 2, \dots, m\}$ . Therefore, the maximum inner product is at least  $\tau$ .

Formulation (D1) has an evident geometric meaning: find the intersection of  $Q$  with the ray  $\{\tau d : \tau \geq 0\}$  by exploring the portion of the line belonging to  $Q$ . Because  $\tau d$  is always required to lie in  $Q$ , this is an *internal* description of the problem. The same underlying geometry can be expressed by considering the portion of the line lying outside  $Q$ , thus arriving at the following *external* description:

$$\boxed{\varphi^* = \min_\tau \{\tau > 0 : \tau d \notin \text{int } Q\}.} \quad (D'1)$$

We mention *both* (D1) and (D'1) because the algorithms we design in later sections give a lower *and* an upper bound on  $\varphi^*$ , thus producing feasible solutions (with certain prescribed

relative accuracy) to *both* problems. In fact, *all* the problems we consider in this paper have optimal value either  $\varphi^*$  or  $1/\varphi^*$  and hence all can be viewed as specific formulations of the same underlying (one-dimensional) geometric problem.

Problem (P1) can be reformulated as follows:

$$\begin{aligned}
\varphi^* &\equiv \min_x \{ \max_i |\langle a_i, x \rangle| \text{ s.t. } \langle d, x \rangle = 1 \} \\
&= \min_{x, \tau} \{ \tau : \max_i |\langle a_i, x \rangle| \leq \tau, \langle d, x \rangle = 1 \} \\
&= \min_{x, \tau} \{ \tau : x \in \tau Q^\circ, \langle d, x \rangle = 1, \tau > 0 \} \\
&= \min_{z, \tau} \{ \tau : z \in Q^\circ, \langle d, z \rangle = 1/\tau, \tau > 0 \} \\
&= \left[ \max_{z, \tau} \{ 1/\tau : z \in Q^\circ, \langle d, z \rangle = 1/\tau, \tau > 0 \} \right]^{-1} \\
&= \left[ \max_z \{ \langle d, z \rangle : z \in Q^\circ \} \right]^{-1},
\end{aligned}$$

and therefore

$$\boxed{1/\varphi^* = \max_z \{ \langle d, z \rangle : z \in Q^\circ \} = \xi_{Q^\circ}(d)}. \tag{P2}$$

If  $x$  is feasible for (P1) then  $z := x/\varphi(x)$  is feasible for (P2) as

$$\max_i |\langle a_i, z \rangle| = \max_i |\langle a_i, x \rangle|/\varphi(x) = 1.$$

On the other hand, if  $z$  is feasible for (P2) then  $x := z/\langle d, z \rangle$  is feasible for (P1) because  $\langle d, x \rangle = \langle d, z \rangle/\langle d, z \rangle = 1$ . A slightly more careful look at the above chain of equalities reveals the following:

**Proposition 1.** *Point  $x = z/\langle d, z \rangle$  is a minimizer of (P1) with optimal value  $\varphi^*$  if and only if  $z = x/\varphi(x)$  is a maximizer of (P2) with optimal value  $1/\varphi^*$ .*

Consider now the dual of (P2). It can be written as

$$\boxed{1/\varphi^* = \min_v \{ \|v\|_1 : Av = d, v \in \mathbf{R}^m \}. \tag{D2}$$

This is the problem of finding the smallest  $\ell_1$ -norm solution of the underdetermined full rank (see assumption (10)) linear system  $Av = d$ . Let us again check weak duality. For any  $z \in Q^\circ$  and  $v$  with  $Av = d$ , one has

$$\|v\|_1 - \langle d, z \rangle = \|v\|_1 - \langle v, A^* z \rangle = \|v\|_1 - \sum_{i=1}^m v^{(i)} \langle a_i, z \rangle \geq \sum_{i=1}^m (|v^{(i)}| - |v^{(i)}| |\langle a_i, z \rangle|) \geq 0.$$

We arrive at this straightforward observation:

**Proposition 2.** *Point  $z$  feasible for (P2) ( $v$  feasible for (D2)) is optimal if and only if there is  $v$  feasible for (D2) ( $z$  feasible for (P2)) such that the following complementary slackness conditions hold:*

$$|v^{(i)}| = v^{(i)} \langle a_i, z \rangle, \quad i = 1, 2, \dots, m. \quad (12)$$

Using the complementary slackness condition (12) between problem (P2) and its dual (D2) together with the relationship between problems (P1) and (P2) given by Proposition 1 and the discussion preceding it, we have arrived at the following complementary slackness condition between problems (P1) and (D2):

$$|v^{(i)}| \varphi(x) = v^{(i)} \langle a_i, x \rangle, \quad i = 1, 2, \dots, m. \quad (13)$$

Note that (13) is equivalent to

$$v^{(i)} \neq 0 \quad \Rightarrow \quad \varphi(x) = |\langle a_i, x \rangle|, \quad \text{and} \quad \text{sign}(\langle a_i, x \rangle) = \text{sign}(v^{(i)}). \quad (14)$$

We have thus shown the following.

**Proposition 3.** *Point  $x$  feasible for (P1) ( $v$  feasible for (D2)) is optimal if and only if there is  $v$  feasible for (D2) ( $x$  feasible for (P1)) such that the following complementary slackness conditions hold for  $i = 1, 2, \dots, m$ :*

$$\begin{aligned} v^{(i)} > 0 &\quad \Rightarrow \quad \varphi(x) = \langle a_i, x \rangle, \quad \text{and} \\ v^{(i)} < 0 &\quad \Rightarrow \quad \varphi(x) = \langle -a_i, x \rangle. \end{aligned}$$

Alternatively, the statement above is equivalent to saying that there is a subdifferential of the objective function  $\varphi$  at  $x$  such that its negative lies in the normal cone to the constraint set at  $x$ .

## 2.2 The main problem

The main focus of this paper is the minimization problem

$$\boxed{\psi^* \stackrel{\text{def}}{=} \min_w \{ \psi(w) \equiv \|d\|_{G(w)}^* : w \in \Delta_m \},} \quad (P3)$$

where

$$G(w) \stackrel{\text{def}}{=} \sum_{i=1}^m w^{(i)} a_i a_i^*. \quad (15)$$

For notational brevity we will henceforth write  $\mathcal{B}(w) = \mathcal{B}(G(w))$  and  $\mathcal{B}^\circ(w) = \mathcal{B}^\circ(G(w))$ .

**Theorem 4.**  $1/\varphi^* = \psi^*$ .

In the remainder of this section we give two different proofs of this claim. Both will give us useful insights. The two simple results that follow merit mentioning due to their usefulness in later sections alone. However, having established them, it takes only a simple minimax argument to prove Theorem 4.

**Fact 5** (Degenerate projection). *If  $G : \mathbf{E} \rightarrow \mathbf{E}^*$  is a self-adjoint semidefinite linear operator, then*

$$\min_{\bar{x}} \{\|\bar{x}\|_G : \langle d, \bar{x} \rangle = 1\} = 0 \quad \Leftrightarrow \quad d \notin \text{range}(G), \quad (16)$$

and the following statements are equivalent:

(i)  $x \in \arg \min_{\bar{x}} \{\|\bar{x}\|_G : \langle d, \bar{x} \rangle = 1\}$ ,  $d \in \text{range}(G)$ ,

(ii)  $Gy = d$ ,  $x = y/(\|d\|_G^*)^2$  for some  $y$ , and

(iii)  $\langle d, x \rangle = \|d\|_G^* \|x\|_G = 1$ .

*Proof.* For example, see [24]. □

**Lemma 6.** *The following hold:*

(i) For all  $x \in \mathbf{E}$  and  $w \in \Delta_m$ ,  $\|x\|_{G(w)} \leq \varphi(x)$ , with equality if and only if the following condition holds:

$$w^{(i)} > 0 \quad \Rightarrow \quad \varphi(x) = |\langle a_i, x \rangle|, \quad i = 1, 2, \dots, m.$$

(ii)  $\varphi(x) = \max_w \{\|x\|_{G(w)} : w \in \Delta_m\}$ .

(iii)  $\mathcal{B}(w) \subset Q$  and  $Q^\circ \subset \mathcal{B}^\circ(w)$  for all  $w \in \Delta_m$ .

*Proof.* Part (i) follows from

$$\xi_{\mathcal{B}(w)}(x) = \|x\|_{G(w)} = \left[ \sum_i w^{(i)} \langle a_i, x \rangle^2 \right]^{1/2} \leq \max_i |\langle a_i, x \rangle| = \xi_Q(x) \equiv \varphi(x).$$

Condition for equality and parts (ii) and (iii) follow easily. □

To prove Theorem 4, observe that

$$\begin{aligned} \varphi^* &= \min_{x: \langle d, x \rangle = 1} \varphi(x) \\ &= \min_{x: \langle d, x \rangle = 1} \max_{w \in \Delta_m} \|x\|_{G(w)} && \text{by part (ii) of Lemma 6} \\ &= \max_{w \in \Delta_m} \min_{x: \langle d, x \rangle = 1} \|x\|_{G(w)} \\ &= \max_{\substack{w \in \Delta_m \\ d \in \text{range } G(w)}} \min_{x: \langle d, x \rangle = 1} \|x\|_{G(w)} && \text{by (16)} \\ &= \max_{\substack{w \in \Delta_m \\ d \in \text{range } G(w)}} 1/\|d\|_{G(w)}^* && \text{by parts (i) and (iii) of Fact 5} \\ &= \max_{w \in \Delta_m} 1/\|d\|_{G(w)}^* \\ &= \left[ \min_{w \in \Delta_m} \|d\|_{G(w)}^* \right]^{-1}. \end{aligned}$$

The interchange of minimum and maximum in the derivation above can be justified using Hartung's [14] generalization of Sion's [29] minimax theorem.

There is another way of seeing that  $\psi^* = 1/\varphi^*$  and that the minimum is attained. The proof reveals a crucial relationship between the feasible solutions of problems (P3) and (D2), suggesting an algorithmic idea for solving both problems. For this we need two intermediate results.

**Lemma 7.** *Assume  $G(w)y = d$  for some  $w \in \Delta_m$  and  $y \in \mathbf{E}$ . If we set  $v^{(i)} = w^{(i)}\langle a_i, y \rangle$  for  $i = 1, \dots, m$ , then  $Av = d$  and  $\|v\|_1 \leq \|d\|_{G(w)}^*$ .*

*Proof.* By the inequality between weighted arithmetic and quadratic means

$$\|v\|_1 = \sum_i w^{(i)} |\langle a_i, y \rangle| \leq \left[ \sum_i w^{(i)} \langle a_i, y \rangle^2 \right]^{1/2} = \langle G(w)y, y \rangle^{1/2} = \langle d, y \rangle^{1/2} = \|d\|_{G(w)}^*.$$

□

This result says: if  $w$  is feasible for (P3) with finite objective value, then the objective value of (D2) for some  $v$  is no bigger than that of (P3) for  $w$ .

**Lemma 8.** *If for  $0 \neq v \in \mathbf{R}^m$  we let  $c := Av$  and  $w := |v|/\|v\|_1$ , then  $\|c\|_{G(w)}^* \leq \|v\|_1$ .*

*Proof.* We can without loss of generality assume that  $\|v\|_1 = 1$  since both sides of the inequality to prove are positively homogeneous in  $v$ . Let  $G = G(|v|)$ . As we will see in Lemma 10, there is  $y \in \mathbf{E}$  for which  $Gy = c$ . The result then follows from

$$\begin{aligned} \langle c, y \rangle &= \left\langle \sum_i v^{(i)} a_i, y \right\rangle \leq \sum_i |v^{(i)}| |\langle a_i, y \rangle| \\ &\leq \left[ \sum_i |v^{(i)}| \langle a_i, y \rangle^2 \right]^{1/2} = \langle Gy, y \rangle^{1/2} = \langle c, y \rangle^{1/2}, \end{aligned}$$

since this implies  $\|c\|_G^* = \langle c, y \rangle^{1/2} \leq 1$ .

□

This result says: if  $v$  is feasible for (D2) then the objective value of (P3) for some  $w$  is no bigger than that of (D2) for  $v$ . The proof of Theorem 4 (and of attainment) is a direct consequence of this for  $c = d$  and the fact that the optimal value of (D2) is  $1/\varphi^*$ .

### 2.3 The seventh problem

If we try to enclose the polytope  $Q$  into an ellipsoidal cylinder (and think of this cylinder as the unit ball with respect to some seminorm on  $\mathbf{E}^*$ ), so that vector  $d$  has as large seminorm as possible, we are trying to solve the following semidefinite program (SDP):

$$\boxed{(\psi^*)^2 = \max\{\langle d, Md \rangle : M \succeq 0; \langle a_i, Ma_i \rangle \leq 1, i = 1, 2, \dots, m\}.}$$
 (D3)

It is intuitively clear that the optimal cylinder will pass through the intersection of  $Q$  and the ray in the direction  $d$ , i.e. it will pass through the point  $\varphi^*d$ . Expressed differently, we expect that  $\langle \varphi^*d, M\varphi^*d \rangle = 1$ , implying that the optimal value of (D3) should indeed be  $1/(\varphi^*)^2 = (\psi^*)^2$ . A rigorous argument and the relation to (P3) can be found in the following result.

**Proposition 9.** *Problem (P3) is equivalent to the Lagrangean dual of (D3) and the optimal value of (D3) is  $(\psi^*)^2$ .*

*Proof.* The first expression in the following chain of identities is the Lagrangean dual of (D3) and the last expression is essentially problem (P3):

$$\begin{aligned} & \min \left\{ \sum_i \bar{w}^{(i)} : \bar{w} \geq 0, \sum_i \bar{w}^{(i)} a_i a_i^* \succeq dd^* \right\} \\ &= \min \{ \tau : \tau > 0, \frac{1}{\tau} dd^* \preceq G(w), w \in \Delta_m \} \\ &\stackrel{\text{(L.14)}}{=} \min \{ \tau : \tau > 0, \tau \geq (\|d\|_{G(w)}^*)^2, w \in \Delta_m \} \\ &= \min \{ (\|d\|_{G(w)}^*)^2 : w \in \Delta_m \}. \end{aligned}$$

In the first step we have passed to the new variable  $w = \bar{w} / \sum \bar{w}^{(i)}$ . □

Also note that if we restrict attention only to rank-one matrices  $M$ , (D3) transforms into (P2) (with an added square in the objective function). This shows that there must be an optimal rank-one solution (which also follows from Theorem 21).

## 2.4 First algorithmic idea

Lemmas 7 and 8 have a nice geometric interpretation: First notice that  $c' = Av/\|v\|_1 \in Q$  and that any point of  $Q \setminus \{0\}$  can be written in this form for some  $0 \neq v \in \mathbf{R}^m$ ; i.e.  $Q = \{0\} \cup \{Av/\|v\|_1 : v \in \mathbf{R}^m \setminus \{0\}\}$ . Also notice that  $0 \in \mathcal{B}(w)$  for all  $w \in \Delta_m$ . Lemma 8 therefore says that any point  $c'$  of  $Q$  can be enclosed into an ellipsoid  $\mathcal{B}(w)$  (by virtue of inequality  $\|c'\|_{G(w)}^* \leq 1$ ) for properly chosen vector of weights  $w$ . To complete the geometric picture, recall that by part (iii) of Lemma 6,  $\mathcal{B}(w) \subset Q$ . In short, any point of  $Q$  is “reachable” by a (possibly degenerate origin-centered) ellipsoid lying in  $Q$ .

## 3 Convexity and smoothness

In this subsection we establish the convexity of the function

$$\psi^2(w) = (\|d\|_{G(w)}^*)^2,$$

and derive formulae for its first and second derivatives. This is the square of the objective function of problem (P3).

### 3.1 Convexity of the domain

Let us start by showing that the domain of  $\psi$  (or equivalently of  $\psi^2$ ), defined the usual way as

$$\text{dom } \psi \stackrel{\text{def}}{=} \{w \in \Delta_m : \psi(w) < +\infty\} = \{w \in \Delta_m : d \in \text{range } G(w)\},$$

is convex. For this we will need two intermediate results.

**Lemma 10.** *For all  $v \in \mathbf{R}^m$ ,  $\text{range } G(v) = \text{span}\{a_i : v^{(i)} \neq 0\}$ .*

*Proof.* Let  $G = G(v)$  and  $\tilde{A}$  be the matrix obtained from  $A = [a_1, \dots, a_m]$  by excluding all columns with zero weights. Let  $\tilde{v}$  be defined in an analogous fashion. Note that for any  $x$ ,  $Gx$  is a linear combination of columns of  $\tilde{A}$  and thus  $\text{range}(G) \subset \text{range}(\tilde{A})$ . However,

$$\text{rank}(G) = \text{rank}(\tilde{A} \text{diag}(\tilde{v}) \tilde{A}^*) = \text{rank}(\tilde{A} \tilde{A}^*) = \text{rank}(\tilde{A})$$

and hence  $\text{range}(G) = \text{range}(\tilde{A})$ . □

Although we have stated the previous result in an ever slightly more general way, we will apply it exclusively in the case  $v = w \in \Delta_m$ . Note that, as a first consequence,  $G(w)$  is invertible (and hence  $\mathcal{B}(w)$  is a full-dimensional ellipsoid) if the vectors  $a_i$  with nonzero weights span  $\mathbf{E}^*$ . Since  $\text{range } A = \mathbf{E}^*$  by (10), this happens, in particular, when all weights are positive.

Let us at this point state a simple corollary to Lemma 10. This is not needed to establish convexity of the domain of  $\psi$ .

**Proposition 11** (Openness of the domain). *The domain of  $\psi$  is open relative to  $\Delta_m$ .*

*Proof.* Let us fix  $w \in \text{dom } \psi$  and note that for all sufficiently small  $h \in \mathbf{R}^m$  we have  $w^{(i)} + h^{(i)} > 0$  whenever  $w^{(i)} > 0$ , for all  $i$ . If, in addition,  $w + h \in \Delta_m$ , then our assumption about  $w$  and Lemma 10 imply  $d \in \text{range } G(w) \subset \text{range } G(w + h)$ . □

**Lemma 12.** *For any  $w_1, w_2 \in \Delta_m$  and  $w = \lambda w_1 + (1 - \lambda)w_2$  with  $0 < \lambda < 1$ ,*

$$\text{range } G(w_1) \cup \text{range } G(w_2) \subset \text{range } G(w).$$

*Proof.* Notice that for any  $i$ , the weight  $w^{(i)}$  is positive if and only if at least one of the weights  $w_1^{(i)}, w_2^{(i)}$  is positive and hence

$$\{a_i : w_1^{(i)} > 0 \text{ or } w_2^{(i)} > 0\} = \{a_i : w^{(i)} > 0\}.$$

By Lemma 10,

$$\begin{aligned} \text{range } G(w_1) \cup \text{range } G(w_2) &= \text{span}\{a_i : w_1^{(i)} > 0\} \cup \text{span}\{a_i : w_2^{(i)} > 0\} \\ &\subset \text{span}\{a_i : w_1^{(i)} > 0 \text{ or } w_2^{(i)} > 0\} \\ &= \text{span}\{a_i : w^{(i)} > 0\} \\ &= \text{range } G(w). \end{aligned}$$

□

**Proposition 13** (Convexity of the domain). *The domain of  $\psi$  is a convex set.*

*Proof.* Assume  $\psi(w_1) < +\infty$  and  $\psi(w_2) < +\infty$  for some  $w_1, w_2 \in \Delta_m$ . This is equivalent to  $d \in \text{range } G(w_1) \cap \text{range } G(w_2)$ . If we consider  $w = \lambda w_1 + (1 - \lambda)w_2$  for  $0 < \lambda < 1$ , then by Lemma 12,  $d \in \text{range } G(w)$  and hence  $\psi(w) < +\infty$ .  $\square$

### 3.2 Convexity of the objective function

**Lemma 14.** *If  $G: \mathbf{E} \rightarrow \mathbf{E}^*$  is positive semidefinite and self-adjoint,  $g \in \mathbf{E}^*$  and  $\tau$  a positive real parameter, then the following statements are equivalent:*

- (i)  $\tau \geq (\|g\|_G^*)^2$ ,
- (ii)  $\begin{pmatrix} \tau & g^* \\ g & G \end{pmatrix} \succeq 0$ ,
- (iii)  $G - \frac{1}{\tau}gg^* \succeq 0$ .

*Proof.* Let us start with the equivalence between (i) and (ii).

If  $g \notin \text{range } G$ , then  $\|g\|_G^* = +\infty$  and we need to show that the operator in (ii) is *not* positive definite for any real  $\tau$ . Since  $G$  is singular,  $\text{null } G$  is nontrivial. Clearly there must be  $y \in \text{null } G$  for which  $\langle g, y \rangle \neq 0$  since  $g \notin \text{range } G = (\text{null } G)^\perp$ . Choose  $y$  with  $\langle g, y \rangle < 0$  and consider

$$\begin{aligned} \begin{pmatrix} \delta & y^* \end{pmatrix} \begin{pmatrix} \tau & g^* \\ g & G \end{pmatrix} \begin{pmatrix} \delta \\ y \end{pmatrix} &= \tau\delta^2 + 2\langle g, y \rangle\delta + \langle Gy, y \rangle \\ &= \tau\delta^2 + 2\langle g, y \rangle\delta, \end{aligned}$$

which is negative for all sufficiently small positive  $\delta$ . Now assume  $g \in \text{range } G$  and take  $y$  such that  $Gy = g$ . For this part of the argument we treat the spaces  $\mathbf{E}$  and  $\mathbf{E}^*$  as  $\mathbf{R}^n$ . We do this in order to be able to use a diagonalization technique. The operator from (ii) is positive semidefinite if and only if the following  $(n + 1) \times (n + 1)$  matrix is positive semidefinite

$$\begin{pmatrix} 1 & -y^T \\ 0 & I_n \end{pmatrix} \begin{pmatrix} \tau & g^T \\ g & G \end{pmatrix} \begin{pmatrix} 1 & 0 \\ -y & I_n \end{pmatrix} = \begin{pmatrix} \tau - \langle g, y \rangle & 0 \\ 0 & G \end{pmatrix}.$$

This happens precisely when  $\tau \geq \langle g, y \rangle = (\|g\|_G^*)^2$ . The equivalence of (ii) and (iii) follows by noting that  $G - \frac{1}{\tau}gg^*$  is the Schur complement of  $\tau$  in the block matrix in (ii).  $\square$

**Proposition 15** (Convexity).  *$\psi^2$  is convex on its domain.*

*Proof.* Consider  $(w_1, \tau_1), (w_2, \tau_2) \in \text{epi } \psi^2$  and  $\lambda \in (0, 1)$ . Letting  $\tau = \lambda\tau_1 + (1 - \lambda)\tau_2$  and  $w = \lambda w_1 + (1 - \lambda)w_2$ , notice that

$$\begin{pmatrix} \tau & d^* \\ d & G(w) \end{pmatrix} = \lambda \begin{pmatrix} \tau_1 & d^* \\ d & G(w_1) \end{pmatrix} + (1 - \lambda) \begin{pmatrix} \tau_2 & d^* \\ d & G(w_2) \end{pmatrix}.$$

Convexity of the epigraph of  $\psi^2$  (and hence of  $\psi^2$ ) now follows from the equivalence of (i) and (ii) in Lemma 14 and the convexity of the cone of positive semidefinite matrices.  $\square$

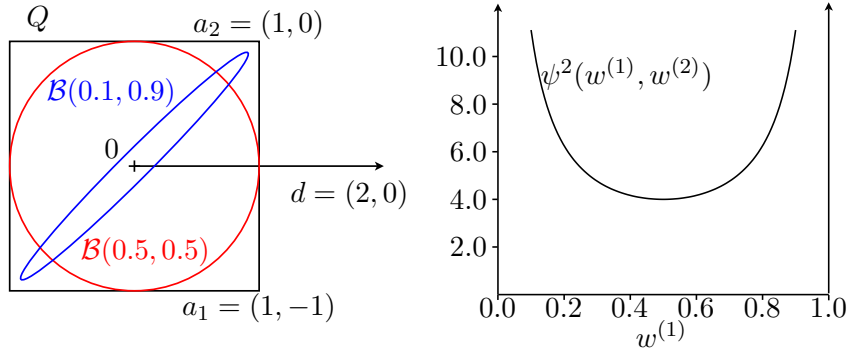


Figure 1: Example 17.

### 3.3 Smoothness

If  $G(w)$  is invertible, then  $\psi^2(w) = \langle d, G(w)^{-1}d \rangle$  is differentiable at  $w$ . For an invertible linear operator  $C: \mathbf{E} \rightarrow \mathbf{E}^*$ , let  $\theta$  be defined by  $\theta(C) = C^{-1}$ . The fact that  $D\theta(C)H = -C^{-1}HC^{-1}$  together with the chain rule gives the following formulae for the first and second (Fréchet) derivatives of  $\psi^2$ .

**Proposition 16** (Differentiability). *If  $G(w)$  is invertible, then  $\psi^2$  is differentiable at  $w$  and for  $h \in \mathbf{R}^m$  we have the following formulae:*

- (i)  $D\psi^2(w)h = -\langle d, G(w)^{-1}G(h)G(w)^{-1}d \rangle$ , and
- (ii)  $D^2\psi^2(w)[h, h] = 2\langle d, G(w)^{-1}G(h)G(w)^{-1}G(h)G(w)^{-1}d \rangle$ .

It is apparent from the form of the Hessian of  $\psi^2$  that it is positive semidefinite. Indeed, for any  $h \in \mathbf{R}^m$  let  $g = G(h)G(w)^{-1}d$  and note that  $D^2\psi^2(w)[h, h] = 2\langle g, G(w)^{-1}g \rangle \geq 0$ . This is an alternative way to establish convexity of  $\psi^2$ , although on a smaller set than  $\text{dom } \psi$ .

### 3.4 An example

**Example 17.** Consider an example with  $n = m = 2$  as in Figure 17. We have  $a_1 = (1, -1)$ ,  $a_2 = (1, 1)$  and  $d = (2, 0)$  and hence for  $(w^{(1)}, w^{(2)}) \in \Delta_2$  we get

$$G(w^{(1)}, w^{(2)}) = w^{(1)}a_1a_1^T + w^{(2)}a_2a_2^T = \begin{pmatrix} 1 & w^{(2)} - w^{(1)} \\ w^{(2)} - w^{(1)} & 1 \end{pmatrix}.$$

Assuming  $w^{(1)} > 0$  and  $w^{(2)} > 0$ , the system  $G(w^{(1)}, w^{(2)})y = d$  has the unique solution

$$y_1 = \frac{1}{2} \left( \frac{1}{w^{(1)}} + \frac{1}{w^{(2)}} \right), \quad y_2 = \frac{1}{2} \left( \frac{1}{w^{(2)}} - \frac{1}{w^{(1)}} \right),$$

and therefore

$$\psi^2(w^{(1)}, w^{(2)}) = \langle d, y \rangle = \frac{1}{w^{(1)}} + \frac{1}{w^{(2)}}.$$

Note that  $\|d\|_{G(0.5,0.5)}^* = \psi(0.5, 0.5) = 2$ , which geometrically corresponds to the ball  $\mathcal{B}(0.5, 0.5)$  cutting vector  $d$  in half (Figure 1). Also observe that as  $w^{(1)} \rightarrow 0$ , the  $\|\cdot\|_{G(w)}^*$ -norm of  $d$  increases to infinity. This translates to the ellipsoid  $\mathcal{B}(w^{(1)}, w^{(2)})$  getting thinner, “approaching” the lower dimensional ellipsoid  $\mathcal{B}(0, 1)$  — the line segment with endpoints  $a_2$  and  $-a_2$ .

If  $w^{(i)} = 0$  then  $\text{range } G(w^{(1)}, w^{(2)}) = \text{span}\{a_{3-i}\}$  and we conclude that in either case  $d \notin \text{range } G(w^{(1)}, w^{(2)})$ , implying  $\psi^2(w^{(1)}, w^{(2)}) = +\infty$ . Notice that  $\psi^2$  is a convex function (as asserted in Proposition 15; also see Figure 1) with convex domain  $\{w \in \Delta_m : w^{(1)} > 0, w^{(2)} > 0\}$  (Proposition 13), which is an open set relative to  $\Delta_m$  (Proposition 11). For any  $(w^{(1)}, w^{(2)}) \in \text{dom } \psi$  and  $h$  with  $h^{(1)} + h^{(2)} = 0$  we have

$$D\psi^2(w^{(1)}, w^{(2)})h = - \left( \frac{h^{(1)}}{(w^{(1)})^2} + \frac{h^{(2)}}{(w^{(2)})^2} \right) = -\langle G(h)y, y \rangle,$$

which agrees with the formula from Proposition 16.

## 4 Optimality conditions

In Section 2 we have outlined the basic relationships among problems  $(P1)$ ,  $(D1)$ ,  $(D'1)$ ,  $(P2)$ ,  $(D2)$ ,  $(P3)$  and  $(D3)$ . In this part we first investigate the necessary and sufficient optimality conditions for problem  $(P3)$  and then show that an approximate version of these implies approximate optimality in the other problems.

**Fact 18** (Cauchy-Schwarz). *For all  $x \in \mathbf{E}$  and  $g \in \text{range}(G)$  we have*

$$\langle g, x \rangle \leq \|g\|_G^* \|x\|_G, \tag{17}$$

with equality exactly in one of the two cases

1.  $\|x\|_G = 0$ , or
2.  $\|x\|_G \neq 0$  and  $g$  is a nonnegative multiple of  $Gx$ .

*Proof.* For example, see [24]. For proof of a more general statement see the definition of a polar gauge and Theorem 15.1 in Rockafellar [26].  $\square$

**Lemma 19.** *Let  $G: \mathbf{E} \rightarrow \mathbf{E}^*$  be self-adjoint and positive semidefinite. Further let  $0 \neq c \in \text{range } G$  and assume that  $y \in \mathbf{E}$  defines a supporting hyperplane to  $\mathcal{B}(G)$  at  $c' := c/\|c\|_G^*$  in the following sense:*

$$\langle g, y \rangle < \langle c', y \rangle \quad \forall c' \neq g \in \mathcal{B}(G). \tag{18}$$

Then  $c = \lambda Gy$  for some  $\lambda > 0$ .

*Proof.* First notice that because  $0 \neq c \in \text{range } G$ , we have  $0 < \|c\|_G^* < \infty$ . Now observe that  $c'$  lies in the relative boundary of  $\mathcal{B}(G)$ , which implies that a vector  $y$  as above exists. By (18)

$$\|y\|_G = \xi_{\mathcal{B}(G)}(y) = \max\{\langle g, y \rangle : g \in \mathcal{B}(G)\} = \langle c', y \rangle,$$

and hence

$$\|y\|_G \|c\|_G^* = \langle c, y \rangle.$$

The condition for equality in the Cauchy-Schwarz inequality (Fact 18) now implies that either  $\|y\|_G = 0$ , or otherwise  $\|y\|_G \neq 0$  and  $c$  is a nonnegative multiple of  $Gy$ . We claim that the first case can be excluded. Indeed, if  $\|y\|_G = 0$  then by (4) we get  $Gy = 0$ , which would in turn imply that  $\langle g, y \rangle = 0$  for all  $g \in \text{range } G \supset \mathcal{B}(G)$ , violating (18). The statement of the lemma then follows from the second case discussed above by noting that the assumption  $c \neq 0$  implies that the nonnegative multiplier must be in fact positive.  $\square$

The above lemma will be used to prove the necessity part of the following optimality condition.

**Theorem 20** (Optimality). *Point  $w \in \Delta_m$  is optimal for (P3) if and only if there exists  $y \in \mathbf{E}$  such that*

$$(i) \quad G(w)y = d, \text{ and}$$

$$(ii) \quad \varphi(y) = \psi(w).$$

Condition (ii) can be replaced by

$$(ii') \quad w^{(i)} > 0 \quad \Rightarrow \quad \varphi(y) = |\langle a_i, y \rangle|, \quad i = 1, 2, \dots, m.$$

*Proof.* If (i) holds then  $x = y/\psi^2(w)$  is feasible for (P1) and hence  $\varphi(y/\psi^2(w)) = \varphi(x) \geq \varphi^* = 1/\psi^*$ . By homogeneity of  $\varphi$  we obtain

$$\varphi(y) \geq \frac{\psi^2(w)}{\psi^*} \geq \psi(w). \tag{19}$$

If we additionally assume (ii) then (19) must hold with equality and thus  $\psi(w) = \psi^*$ . As a side product, this also shows that  $x$  is optimal for (P1). Conversely, assume  $w \in \Delta_m$  is a minimizer of (P3). Then  $\|d\|_{G(w)}^* = \psi^* = 1/\varphi^*$  and  $\varphi^*d \in \text{bdry } Q$  by (11). Let  $y \in \mathbf{E}$  define a supporting hyperplane to  $Q$  at  $c' := \varphi^*d$  (and hence to  $\mathcal{B}(G(w))$  at the same point) so that  $\langle \cdot, y \rangle$  is maximized over  $Q$  at  $c'$  (and hence uniquely over  $\mathcal{B}(G(w))$  at the same point). Applying Lemma 19 with  $G := G(w)$  and  $c := d$  we conclude that  $d = \lambda G(w)y$  for some positive  $\lambda$ . Let us scale  $y$  so that  $d = G(w)y$ , establishing (i). Part (ii) follows from

$$\varphi(y) = \xi_Q(y) = \max\{\langle g, y \rangle : g \in Q\} = \langle c', y \rangle = \langle \varphi^*d, y \rangle = \varphi^*\psi^2(w) = \psi(w).$$

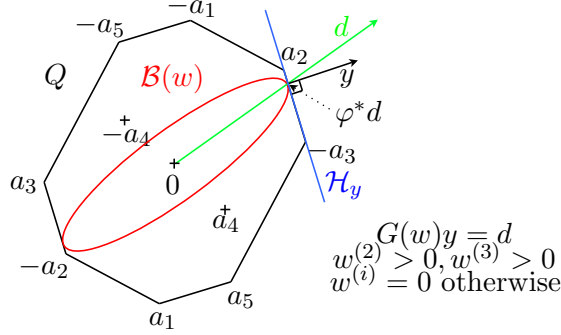


Figure 2: Geometry at optimality.

The first equality is the definition of  $\varphi$ , the third is a consequence of the choice of  $y$  and the last one follows from  $\psi(w) = \psi^* = 1/\varphi^*$ . Finally, the equivalence of (ii) and (ii'), assuming (i), is apparent if we note that  $\varphi(y) = \max_i \{ |\langle a_i, y \rangle| : i = 1, 2, \dots, m \}$  and

$$\psi(w) = \|d\|_{G(w)}^* = \|y\|_{G(w)} = \left( \sum_i w^{(i)} \langle a_i, y \rangle^2 \right)^{1/2}.$$

(See also part (i) of Lemma 6). □

The optimality conditions of the previous result have a clear geometric meaning (see Figure 2). A point  $w \in \Delta_m$  is optimal for (P3) precisely when there exists a hyperplane  $\mathcal{H}_y$  passing through  $d/\|d\|_{G(w)}^*$  which also happens to be a supporting hyperplane of  $Q$ . The set  $\mathcal{F}_y := \mathcal{H}_y \cap Q$  is therefore a face of  $Q$  exposed by the direction  $y$ . Optimality condition (ii') then requires one of the points  $a_i$  or  $-a_i$  to lie in  $\mathcal{F}_y$  if the corresponding weight  $w^{(i)}$  is positive. In other words, if both  $a_i$  and  $-a_i$  lie outside this face, then they must have zero weights, at optimality.

Note also that for optimal  $w$ , the point  $v \in \mathbf{R}^m$  defined by  $v^{(i)} = w^{(i)} \langle a_i, y \rangle$  is optimal for (D2), which is a consequence of Lemma 7 and the fact that the optimal values of problems (P3) and (D2) are equal. The intersection point  $\varphi^* d$  of  $\{\tau d : \tau \geq 0\}$  and  $Q$  can be written as

$$\varphi^* d = \varphi^* G(w)y = \varphi^* \sum_{i=1}^m w^{(i)} \langle a_i, y \rangle a_i = \sum_{i=1}^m \varphi^* v^{(i)} a_i$$

with  $\|\varphi^* v\|_1 = \varphi^* \|v\|_1 = \varphi^* \psi^* = 1$ . Hence the point  $\varphi^* d$  can be written as a convex combination of points  $\pm a_i$  lying in  $\mathcal{F}_y$ , implying that it also lies on the face.

Our next result states that once we are in the possession of  $w$  satisfying condition (i) of Theorem 20 and *approximately* satisfying condition (ii) therein, then  $w$  is *approximately* optimal for (P3) and it is easy to construct nearly optimal solutions to the remaining five problem (P1), (D1), (D'1), (P2) and (D2). In this sense, this is a (much) stronger version of the sufficiency part of Theorem 20.

**Theorem 21** (Approximate optimality). *Let  $G(w)y = d$  for some  $w \in \Delta_m$  and  $y \in \mathbf{E}$  and assume the following  $\delta$ -approximate optimality condition holds:*

$$\varphi(y) \leq (1 + \delta)\psi(w). \quad (20)$$

*Then the points  $x, \tau, \tau', z, v, w$  and  $M$  defined by*

$$(i) \ x = y/\psi^2(w) \quad (P1),$$

$$(ii) \ \tau = 1/\psi(w) \quad (D1),$$

$$(ii') \ \tau' = \varphi(y)/\psi^2(w) = \varphi(x) \quad (D'1),$$

$$(iii) \ z = y/\varphi(y) \quad (P2),$$

$$(iv) \ v \in \mathbf{R}^m \text{ given by } v^{(i)} = w^{(i)}\langle a_i, y \rangle \quad (D2),$$

$$(v) \ w \quad (P3), \text{ and}$$

$$(vi) \ M = yy^*/\psi^2(w) \quad (D3)$$

*are feasible in their respective problems, with the exception of  $M$  which is approximately feasible in the following sense:*

$$M \succeq 0, \quad \langle a_i, Ma_i \rangle \leq (1 + \delta)^2, \quad i = 1, \dots, m, \quad (21)$$

*and satisfy the following  $\delta$ -optimality conditions:*

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

$$(ii) \ \tau \geq (1 + \delta)^{-1}\varphi^* \geq (1 - \delta)\varphi^*,$$

$$(ii') \ \tau' \leq (1 + \delta)\varphi^*,$$

$$(iii) \ \langle d, z \rangle \geq (1 + \delta)^{-1}\psi^* \geq (1 - \delta)\psi^*,$$

$$(iv) \ \|v\|_1 \leq (1 + \delta)\psi^*,$$

$$(v) \ \psi(w) \leq (1 + \delta)\psi^*, \text{ and}$$

$$(vi) \ (\psi^*)^2 \leq \langle d, Md \rangle \leq (1 + \delta)^2(\psi^*)^2.$$

*Moreover, the seven objective values are related in the following way:*

$$\|v\|_1 \leq \frac{1}{\tau} = \psi(w) = \sqrt{\langle d, Md \rangle} \leq (1 + \delta)\frac{1}{\varphi(x)} = (1 + \delta)\frac{1}{\tau'} = (1 + \delta)\langle d, z \rangle.$$

*Proof.* Feasibility in each case follows from the definition of the corresponding point. The non-evident part of the approximate feasibility relations (21) is established as follows:

$$\max_i \langle a_i, Ma_i \rangle = \frac{\max_i \langle a_i, y \rangle^2}{\psi^2(w)} = \frac{\varphi^2(y)}{\psi^2(w)} \stackrel{(20)}{\leq} (1 + \delta)^2.$$

Since  $x = y/\psi^2(w)$  and  $\psi(w) \geq \psi^* = 1/\varphi^*$ , condition (20) yields

$$\varphi(x) \leq (1 + \delta) \frac{1}{\psi(w)} \leq (1 + \delta) \varphi^*,$$

establishing (i). Part (ii)' then follows from (i) as  $\tau' = \varphi(x)$ . Reversing the first of the displayed inequalities above gives (v) since  $\varphi(x) \geq \varphi^*$ . By definition,  $\tau = 1/\psi(w)$  and hence (ii) can be obtained from (v) by taking the reciprocals and substituting  $\psi^* = 1/\varphi^*$ . Part (iii) follows from (i) by noting that  $z = y/\varphi(y)$ ,  $x = y/\psi^2(w)$  and  $\langle d, x \rangle = 1$  implies

$$\langle d, z \rangle = \frac{\langle d, y \rangle}{\varphi(y)} = \frac{\langle d, x \rangle \psi^2(w)}{\varphi(x) \psi^2(w)} = \frac{1}{\varphi(x)}.$$

Inequality (iv) follows from (v) and Lemma 7. To establish (vi), note that  $\langle d, Md \rangle = \langle d, y \rangle = \psi^2(w)$ . It remains to apply feasibility of  $w$  and inequality (v). The final statement can be easily extracted from the proof.  $\square$

Observe that whenever  $\varphi(y) = \psi(w)$ , the values of  $\tau$  and  $\tau'$  are equal. The reason for defining the former as  $1/\psi(w)$  and the latter as  $\varphi(x)$  is because for any  $w$  feasible for (P3) the value  $1/\psi(w)$  gives a lower bound on  $\varphi^*$ , thus producing a feasible point for (D1), while for any  $x$  feasible for (P1) the value  $\varphi(x)$  always yields an upper bound on  $\varphi^*$ , giving a feasible point for (D'1).

## 5 Rank-one update

### 5.1 A multiplicative weight update algorithm

The inequalities formulated as Lemma 7 and Lemma 8 reveal a close relationship between the feasible solutions of problems (D2) and (P3). As we have seen, these two lemmas can be used to argue that the two problems have equal optima (see Theorem 4). However, they are more interesting to us because of their geometry and algorithmic implications.

Assuming we start with a feasible solution to problem (P3), Lemma 7 provides us with a feasible solution to problem (D2) with a better objective value. Now in turn, starting from this feasible solution, Lemma 8 gets us back to a feasible solution to problem (P3), again with a better objective value. Using this observation we have arrived at our first algorithm (Algorithm 1), updating the weights in a multiplicative fashion at every iteration.

Due to their simplicity, multiplicative weight update algorithms have been proposed in the literature for many computer science problems. For a recent unifying review of such approaches we refer the reader to [2].

---

**Algorithm 1 (MultWeight)** Multiplicative weight updates

---

- 1: **Input:**  $a_1, \dots, a_m \in \mathbf{E}^*$ ,  $d \in \mathbf{E}^*$ ,  $\delta > 0$ ;
  - 2: **Initialize:**  $k = 0$ ,  $w_0 = e/m$ ;
  - 3: **Iterate:**
  - 4:  $G_k = \sum_i w_k^{(i)} a_i a_i^*$ ,  $y_k = G_k^{-1} d$ ;
  - 5:  $\alpha_k = \langle d, y_k \rangle$ ,  $j = \arg \max_i |\langle a_i, y_k \rangle|$ ,  $\beta_k = \langle a_j, y_k \rangle$ ;
  - 6:  $\delta_k = \frac{|\beta_k|}{\sqrt{\alpha_k}} - 1$ ;
  - 7: **if**  $\delta_k \leq \delta$
  - 8:     **terminate**;
  - 9: **else**
  - 10:      $v^{(i)} = w_k^{(i)} \langle a_i, y_k \rangle$ ,  $i = 1, 2, \dots, m$ ;
  - 11:      $w_{k+1} = |v| / \|v\|_1$ ;
  - 12:      $k \leftarrow k + 1$ ;
  - 13: **end if**
  - 14: **Output:**  $w_k$  satisfying  $\|d\|_{G(w_k)}^* = \sqrt{\alpha_k} \leq (1 + \delta)\psi^*$ ;
- 

Notice that the stopping criterion of Algorithm 1 is equivalent to the condition (20) of Theorem 21 and hence the point  $w_k$  output by the algorithm, if it terminates, is a  $\delta$ -approximate minimizer of problem (P3). The algorithm, however, suffers from at least two shortcomings.

First, it can fail to terminate. In short, this is because once a weight is set to zero, it can never be increased to a nonzero value, even if the corresponding point  $a_i$  is required to have a positive weight in the optimum. Imagine we run the algorithm starting with positive weights only for  $i \in \mathcal{I}$  with  $\mathcal{I}$  being a proper subset of the index set  $\{1, 2, \dots, m\}$ . It is clear that the algorithm will never be able to work with points  $a_i$  with  $i \notin \mathcal{I}$  and hence we are actually at best trying to find the intersection of the half-line emanating from the origin in the direction  $d$  and the convex hull of  $\{\pm a_i, i \in \mathcal{I}\}$ , which is a proper subset of  $Q$ . If the algorithm happens to drop weights to zero along the way, it will never be able to recover them back to a nonzero value. Because  $\mathcal{I}_{k+1} \subseteq \mathcal{I}_k$  holds for all  $k$ , the “scope” of the method will be the gradually diminishing convex set  $Q_{\mathcal{I}_k} = \text{conv}\{\pm a_i : i \in \mathcal{I}_k\}$ .

Another obvious disadvantage of the algorithm is its high computational cost per iteration due to the need to update  $G$  in a full-rank fashion. The inverse (or a factorization) of  $G$  will therefore have to be fully recomputed at every iteration at the cost of at least  $O(n^3)$  arithmetic operations.

## 5.2 Ingredients of a rank-one update algorithm

As we have already mentioned above, the multiplicative weight update algorithm has the obvious disadvantage of altering weights in a rather *nonuniform* way, resulting in the need to *fully* resolve a system of the form  $G(w)y = d$  at every iteration. The idea we are going to exploit now is updating  $G(w)$  only slightly at every iteration, in a rank-one fashion. This corresponds to changing the weight of a specific term  $a_j a_j^*$  and then adjusting all other weights *uniformly* by a

certain factor, so as to keep the resulting vector of weights feasible.

In what follows we will focus on a single iteration with current weight  $w$ . Assume throughout that  $\psi(w) < \infty$ , or equivalently,

$$d \in \text{range } G(w), \quad (\text{Assumption 1})$$

and that we are in possession of vector  $y$  such that  $G(w)y = d$ . Suppose we update this weight to

$$w(\kappa) \stackrel{\text{def}}{=} \frac{w + \kappa e_j}{1 + \kappa}, \quad (22)$$

where  $\kappa$  is a real parameter,  $j \in \{1, 2, \dots, m\}$  is to be determined later and  $e_j$  is the  $j$ -th unit vector of  $\mathbf{R}^m$ . The smallest possible  $\kappa$  for which  $w(\kappa)$  is feasible is  $\kappa_{\min} := -w^{(j)}$ . For  $w(\kappa)$  to be meaningfully defined, we will further suppose that

$$w^{(j)} \neq 1. \quad (\text{Assumption 2})$$

This ensures both that  $w(\kappa)$  varies as  $\kappa$  varies and that  $w(-w^{(j)})$  is well-defined. We allow  $\kappa$  to take on the value  $\infty$  and naturally define  $w(\infty) := e_j$ . Note that the set of weights described this way forms a chord of  $\Delta_m$  joining vertex  $e_j$  with  $w$ . We chose this particular parametrization of the chord over the more natural  $w(\lambda) := (1 - \lambda)w + \lambda e_j$ ,  $0 \leq \lambda \leq 1$ , because it turns out to yield a more compact exact line-search formula, developed in the next subsection. By linearity of  $G(\cdot)$  as a function of  $w$ , this translates into updating  $G(w)$  as follows:

$$G(\kappa) \stackrel{\text{def}}{=} G(w(\kappa)) = \frac{G(w) + \kappa a_j a_j^*}{1 + \kappa}. \quad (23)$$

This notation, we hope, should not lead to be confusion as we will always use either  $w$  or  $\kappa$  to denote which object do we have in mind.

After we update  $w$  to  $w(\kappa)$ , the value  $\psi^2(w) = \langle d, y \rangle$  changes to

$$\psi^2(\kappa) \stackrel{\text{def}}{=} \psi^2(w(\kappa)) = \begin{cases} \langle d, y(\kappa) \rangle & \text{if } y(\kappa) \text{ solves } G(\kappa)y(\kappa) = d, \\ \infty & \text{if } d \notin \text{range } G(\kappa). \end{cases}$$

Of course, we consider only updates decreasing the objective value at every iteration and hence the second case does not apply for  $\kappa$  we actually end up using.

Since  $G(w)$  changes in a simple and highly structured way (rank-one update and scaling), it can be expected that  $y(\kappa)$  should be obtainable from  $y$  with less effort than resolving from scratch. This is indeed the case. If both  $G(w)$  and  $G(\kappa)$  are nonsingular, one can use the Sherman-Morrison formula (see, for example, Section 2.1.3 of [12]) to this purpose. Loosely speaking, the formula says that the inverse of a rank-one perturbation of a nonsingular matrix results in a rank-one perturbation of the inverse:

**Fact 22** (Sherman-Morrison). *If  $G: \mathbf{E} \rightarrow \mathbf{E}^*$  is an invertible linear operator,  $g \in \mathbf{E}^*$  and  $\kappa \in \mathbf{R}$  such that  $1 + \kappa \langle g, G^{-1}g \rangle \neq 0$ , then  $G + \kappa g g^*$  is invertible and*

$$(G + \kappa g g^*)^{-1} = G^{-1} - \frac{\kappa G^{-1} g g^* G^{-1}}{1 + \kappa \langle g, G^{-1}g \rangle}.$$

In the definition of  $G(\kappa)$  we are dealing with a rank-one update followed by scaling. In particular, if  $G = G(w)$  is invertible and  $1 + \kappa\langle a_j, G^{-1}a_j \rangle \neq 0$  with  $\kappa$  being a real number, the Sherman-Morrison formula implies

$$y(\kappa) := G(\kappa)^{-1}d = (1 + \kappa) \left( y - \frac{\kappa G^{-1}a_j \langle a_j, y \rangle}{1 + \kappa \langle a_j, G^{-1}a_j \rangle} \right),$$

and hence

$$\psi^2(\kappa) = \langle d, y(\kappa) \rangle = (1 + \kappa) \left( \langle d, y \rangle - \frac{\kappa \langle a_j, y \rangle^2}{1 + \kappa \langle a_j, G^{-1}a_j \rangle} \right). \quad (24)$$

In the remainder of this subsection we compute a general formula for  $\psi^2(\kappa)$ , one that is free of the full-rank assumption on  $G$  and includes the case  $\kappa = \infty$  and the situation when the expression in the denominator of (24) vanishes. We proceed through several auxiliary results — the first step is the following simple generalization of the Sherman-Morrison inversion identity:

**Lemma 23.** *Let  $G: \mathbf{E} \rightarrow \mathbf{E}^*$  be a (not necessarily invertible) linear operator and assume  $Gy = d$  for some  $y \in \mathbf{E}$  and  $d \in \mathbf{E}^*$ . If for  $g \in \mathbf{E}^*$  and  $\kappa \in \mathbf{R}$  we let*

$$\tilde{y}(\kappa) := \begin{cases} y & \text{if } \langle g, y \rangle = 0, \\ y - \frac{\kappa \langle g, y \rangle x}{1 + \kappa \langle g, x \rangle} & \text{if } Gx = g \text{ and } 1 + \kappa \langle g, x \rangle \neq 0 \text{ for some } x \in \mathbf{E}, \end{cases}$$

then  $(G + \kappa gg^*)\tilde{y}(\kappa) = d$ .

*Proof.* The first case is trivial; the statement in the second case follows from:

$$\begin{aligned} (G + \kappa gg^*)\tilde{y}(\kappa) &= Gy + \kappa \langle g, y \rangle g - \frac{\kappa \langle g, y \rangle Gx}{1 + \kappa \langle g, x \rangle} - \frac{\kappa^2 \langle g, x \rangle \langle g, y \rangle g}{1 + \kappa \langle g, x \rangle} \\ &= d + \kappa \langle g, y \rangle g \left( 1 - \frac{1 + \kappa \langle g, x \rangle}{1 + \kappa \langle g, x \rangle} \right) = d. \end{aligned}$$

□

**Remark 24.** *Note that if  $G$  is self-adjoint, the value  $\langle g, x \rangle$  in the above lemma does not depend on the particular choice of the solution of the system  $Gx = g$ . Indeed, if  $x'$  and  $x''$  are two such solutions, then  $\langle g, x' \rangle = \langle Gx'', x' \rangle = \langle Gx', x'' \rangle = \langle g, x'' \rangle$ . This is precisely one of the two arguments we used to show that (5) gives a valid definition of  $\|g\|_G^*$ . If we also have  $G \succeq 0$ , then  $\langle g, x \rangle = (\|g\|_G^*)^2$ , which is positive unless  $g = 0$ .*

The next result characterizes the family of rank-one self-adjoint perturbations of a positive semidefinite self-adjoint operator preserving positive-semidefiniteness.

**Lemma 25.** *Let  $G: \mathbf{E} \rightarrow \mathbf{E}^*$  be a positive semidefinite self-adjoint operator and consider  $g \in \mathbf{E}^*$  and a real parameter  $\kappa$ .*

(i) *If  $g \in \text{range } G$ , then*

$$G + \kappa gg^* \succeq 0 \quad \Leftrightarrow \quad 1 + \kappa (\|g\|_G^*)^2 \geq 0.$$

(ii) If  $g \notin \text{range } G$ , then

$$G + \kappa gg^* \succeq 0 \quad \Leftrightarrow \quad \kappa \geq 0.$$

*Proof.* The statements trivially hold if  $\kappa \geq 0$ . If we notice that  $\|g\|_G^* = \infty$  precisely when  $g \notin \text{range } G$ , the case with  $\kappa < 0$  is essentially a restatement of the equivalence between (i) and (iii) of Lemma 14 with  $\tau := -1/\kappa > 0$ .  $\square$

**Corollary 26.** Assume  $a_j \in \text{range } G(w)$ . If  $w(\kappa)$  is feasible, then  $\kappa \geq -1/(\|a_j\|_{G(w)}^*)^2$  and in particular  $-w^{(j)} \geq -1/(\|a_j\|_{G(w)}^*)^2$ . If  $w^{(j)} > 0$ , then  $(\|a_j\|_{G(w)}^*)^2 \leq 1/w^{(j)}$ .

*Proof.* Observe that  $w(\kappa)$  being feasible implies  $1 + \kappa \geq 1 - w^{(j)} > 0$  and  $G(\kappa) \succeq 0$  and hence  $G + \kappa a_j a_j^* = (1 + \kappa)G(\kappa) \succeq 0$ . Now use Lemma 25 with  $g = a_j$ .  $\square$

**Proposition 27.** Let  $G: \mathbf{E} \rightarrow \mathbf{E}^*$  be a positive semidefinite self-adjoint operator and assume  $Gy = d$  for some  $y \in \mathbf{E}$  and  $0 \neq d \in \mathbf{E}^*$ .

(i) If  $0 \neq g \in \text{range } G$  then for  $\kappa \geq -1/(\|g\|_G^*)^2$  the operator  $G + \kappa gg^*$  is positive semidefinite and

$$\|d\|_{G+\kappa gg^*}^*{}^2 = \begin{cases} (\|d\|_G^*)^2 - \frac{\kappa \langle g, y \rangle^2}{1 + \kappa (\|g\|_G^*)^2} & \text{if } \kappa > -\frac{1}{(\|g\|_G^*)^2}, \\ (\|d\|_G^*)^2 & \text{if } \kappa = -\frac{1}{(\|g\|_G^*)^2}, \langle g, y \rangle = 0, \\ \infty & \text{if } \kappa = -\frac{1}{(\|g\|_G^*)^2}, \langle g, y \rangle \neq 0. \end{cases} \quad (25)$$

Moreover,

$$\|d\|_{G+\kappa gg^*}^*{}^2 \rightarrow \begin{cases} \frac{(\|d\|_G^*)^2 (\|g\|_G^*)^2 - \langle g, y \rangle^2}{(\|g\|_G^*)^2} & \text{as } \kappa \rightarrow \infty, \\ \infty & \text{as } \kappa \downarrow -\frac{1}{(\|g\|_G^*)^2} \quad \text{if } \langle g, y \rangle \neq 0. \end{cases}$$

(ii) If  $g \notin \text{range } G$ , then for  $\kappa \geq 0$  the operator  $G + \kappa gg^*$  is positive semidefinite and

$$\|d\|_{G+\kappa gg^*}^* = \|d\|_G^*. \quad (26)$$

*Proof.* Consider statement (i) and note that  $g \neq 0$  implies  $\|g\|_G^* > 0$ . Positive semidefiniteness of  $G + \kappa gg^*$  follows from part (i) of Lemma 25. If  $\tilde{y}(\kappa)$  and  $x$  are as in Lemma 23, then  $(\|d\|_{G+\kappa gg^*}^*)^2 = \langle d, \tilde{y}(\kappa) \rangle$  and  $(\|g\|_G^*)^2 = \langle g, x \rangle$ , and hence the first two cases of (25) follow. Assume now that  $\kappa = -1/(\|g\|_G^*)^2$  and  $\langle g, y \rangle \neq 0$ . We will show that this implies  $d \notin \text{range}(G + \kappa gg^*)$ , and hence  $\|d\|_{G+\kappa gg^*}^* = \infty$ , by demonstrating that  $x' := x/\langle g, y \rangle$  satisfies  $\langle d, x' \rangle = 1$  and  $\langle (G + \kappa gg^*)x', x' \rangle = 0$  and then appealing to Fact 5. Indeed,

$$\langle d, x' \rangle = \frac{\langle d, x \rangle}{\langle g, y \rangle} = \frac{\langle Gy, x \rangle}{\langle g, y \rangle} = \frac{\langle Gx, y \rangle}{\langle g, y \rangle} = 1$$

and

$$\langle (G + \kappa gg^*)x', x' \rangle = \frac{\langle g, x \rangle + \kappa \langle g, x \rangle^2}{\langle g, y \rangle^2} = \frac{(\|g\|_G^*)^2}{\langle g, y \rangle^2} (1 + \kappa (\|g\|_G^*)^2) = 0.$$

The proof of the limit statements is straightforward. To establish (ii), fix arbitrary nonnegative  $\kappa$  and note that whenever some  $\tilde{y}(\kappa)$  satisfies  $(G + \kappa gg^*)\tilde{y}(\kappa) = d$ , we have

$$\kappa \langle g, \tilde{y}(\kappa) \rangle g = d - G\tilde{y}(\kappa) \in \text{range } G.$$

This is possible if and only if  $\kappa \langle g, \tilde{y}(\kappa) \rangle = 0$  and  $G\tilde{y}(\kappa) = d$ . It therefore follows that  $\|d\|_{G+\kappa gg^*}^* = \langle d, \tilde{y}(\kappa) \rangle^{1/2} = \|d\|_G^*$ .  $\square$

The main result of this subsection gives a complete characterization of  $\psi^2(\kappa)$ , generalizing (24).

**Theorem 28.** *Under Assumptions 1 and 2, let  $y \in \mathbf{E}$  be such that  $G(w)y = d$  and let us establish the following simplified notation:*

$$\alpha := \langle d, y \rangle = (\|d\|_{G(w)}^*)^2 = \psi^2(w), \quad \beta := \langle a_j, y \rangle, \quad \gamma := (\|a_j\|_{G(w)}^*)^2. \quad (27)$$

(i) *If  $a_j \in \text{range } G(w)$  and  $-\frac{1}{\gamma} \leq \kappa \leq \infty$ , then the operator  $G(\kappa)$  defined in (23) is positive semidefinite and self-adjoint and  $\psi^2(\kappa) = (\|d\|_{G(\kappa)}^*)^2$  can be written explicitly in terms of  $\alpha, \beta, \gamma$  and  $\kappa$  as follows:*

$$\psi^2(\kappa) = \begin{cases} (1 + \kappa) \left( \alpha - \frac{\kappa \beta^2}{1 + \kappa \gamma} \right) & \text{if } \infty > \kappa > -\frac{1}{\gamma}, \\ (1 + \kappa) \alpha & \text{if } \kappa = -\frac{1}{\gamma}, \beta = 0, \\ \infty & \text{if } \kappa = -\frac{1}{\gamma}, \beta \neq 0, \\ \infty & \text{if } \kappa = \infty, \alpha \gamma > \beta^2, \\ \frac{\alpha}{\gamma} & \text{if } \kappa = \infty, \alpha \gamma = \beta^2. \end{cases} \quad (28)$$

Moreover,  $\psi^2$  enjoys the following continuity/barrier properties:

$$\psi^2(\kappa) \rightarrow \begin{cases} \infty & \text{as } \kappa \downarrow -\frac{1}{\gamma} \text{ if } \beta \neq 0, \\ \infty & \text{as } \kappa \rightarrow \infty \text{ if } \alpha \gamma > \beta^2, \\ \frac{\alpha}{\gamma} & \text{as } \kappa \rightarrow \infty \text{ if } \alpha \gamma = \beta^2. \end{cases} \quad (29)$$

(ii) *If  $a_j \notin \text{range } G(w)$ , and  $0 \leq \kappa \leq \infty$ , then the operator  $G(\kappa)$  is positive semidefinite and self-adjoint and  $\psi^2(\kappa)$  can be written as follows:*

$$\psi^2(\kappa) = \begin{cases} (1 + \kappa) \alpha & \text{if } \infty > \kappa \geq 0, \\ \infty & \text{if } \kappa = \infty. \end{cases} \quad (30)$$

Moreover,  $\psi^2(\kappa) \rightarrow \infty$  as  $\kappa \rightarrow \infty$ .

*Proof.* The first three cases of (28) follow from Proposition 27 used with  $G = G(w)$  and  $g = a_j$  since

$$\psi^2(\kappa) = (\|d\|_{G(\kappa)}^*)^2 = (1 + \kappa)(\|d\|_{G(w) + \kappa a_j a_j^*}^*)^2.$$

The first limit case of (29) corresponds to a case from Proposition 27 while the other two can be easily derived by taking the limit in the first expression of (28).

It remains to analyze the  $\kappa = \infty$  cases. First observe that  $G(\infty) = a_j a_j^*$  and that by the Cauchy-Schwarz inequality (Fact 18)

$$\beta^2 = \langle a_j, y \rangle^2 \leq (\|a_j\|_G^*)^2 \|y\|_G^2 = \gamma \alpha, \quad (31)$$

with equality if and only if either  $a_j$  or  $-a_j$  is a nonnegative multiple of  $Gy = d$  (i.e.  $a_j$  and  $d$  are collinear). Consider the equality case and assume  $d = \tau a_j$ . Since the ellipsoid  $\mathcal{B}(a_j a_j^*)$  corresponds to the line segment  $[-a_j, a_j]$ , it must be the case that  $\|d\|_{a_j a_j^*}^* = |\tau|$ . This can be also seen without referring to the geometrical picture as follows: If  $y'$  is such that  $(a_j a_j^*)y' = d$ , then  $\tau = \langle a_j, y' \rangle$  and

$$(\|d\|_{a_j a_j^*}^*)^2 = \langle d, y' \rangle = \langle \tau a_j, y' \rangle = \tau^2.$$

If  $\langle a_j, y \rangle^2 = \beta^2 = \alpha \gamma > 0$ , we can write

$$\|d\|_{G(\infty)}^* = \tau = \left| \frac{\tau \langle a_j, y \rangle}{\langle a_j, y \rangle} \right| = \frac{\langle d, y \rangle}{|\langle a_j, y \rangle|} = \frac{\alpha}{|\beta|},$$

and hence

$$\psi^2(\infty) = (\|d\|_{G(\infty)}^*)^2 = \frac{\alpha^2}{\beta^2} = \frac{\alpha^2}{\alpha \gamma} = \frac{\alpha}{\gamma}.$$

In the remaining case  $a_j$  and  $d$  are not collinear and thus  $d \notin \text{range}(G(\infty))$ , implying that  $\psi^2(\infty) = \infty$ .

The first statement of part (ii) is a consequence of part (ii) of Proposition 27. The second statement can be proved in complete analogy to the fourth case of (28). This is because  $d \in \text{range } G$  and  $a_j \notin \text{range } G$  and hence  $d$  and  $a_j$  can not be collinear, implying  $\alpha \gamma > \beta^2$ .  $\square$

## 6 Line search

In this subsection we consider the following line-search problem:

$$\boxed{\kappa^* := \arg \min \{ \psi^2(\kappa) : \kappa_{min} := -w^{(j)} \leq \kappa \leq \infty \}.} \quad (32)$$

As in the previous section, we will adhere to Assumptions 1 and 2. Note that by Corollary 26,  $-1/\gamma \leq -w^{(j)} = \kappa_{min}$ , and hence for the purposes of the line-search problem above,  $\psi^2(\kappa)$  is fully described by Theorem 28. First observe that if  $a_j \notin \text{range } G(w)$  then  $w^{(j)} = 0$  by Lemma 10. In this case, in view of Theorem 28, the line-search problem is trivial with the optimal step size being  $\kappa^* = 0$ . We will thus henceforth in this section assume that

$$a_j \in \text{range } G(w). \quad (\text{Assumption 3})$$

Our main result in Subsection 6.1 is Theorem 29, in which we give a closed-form formula for the solution of (32), for arbitrary  $j$ . In the next subsection we then specialize this formula for two important choices of  $j$  on which we eventually base our algorithms.

## 6.1 General line-search formula

For simplicity of the analysis that follows, we will use notation introduced in (27). Note that  $0 < \alpha < +\infty$  (the former inequality is due to  $d \neq 0$  and the latter because of Assumption 1) and  $0 < \gamma < +\infty$  (the former due to Assumption 3 and the latter because  $a_j \neq 0$ ). In fact, recall that we assume throughout the paper that *all* the vectors  $a_1, \dots, a_m$  are nonzero. Also recall that by the Cauchy-Schwarz inequality,  $\alpha\gamma \geq \beta^2$  (31), with equality if and only if  $a_j$  and  $d$  are collinear.

**Theorem 29.** *Under Assumptions 1,2 and 3, the solution of the line-search problem (32) is*

$$\kappa^* = \begin{cases} \kappa_{min} & \text{if } \beta = 0 \text{ or } \gamma \leq 1, \text{ and otherwise} \\ \max\{\kappa_{min}, \kappa_1\} & \text{if } \alpha\gamma > \beta^2, \\ \infty & \text{if } \alpha\gamma = \beta^2, \end{cases} \quad (33)$$

where

$$\kappa_1 := -\frac{1}{\gamma} + \frac{|\beta|\sqrt{\gamma-1}}{\gamma\sqrt{\alpha\gamma-\beta^2}}. \quad (34)$$

Moreover, if  $-1/\gamma = -w^{(j)}$  then  $\gamma > 1$ . If, additionally,  $\alpha\gamma > \beta^2$ , then  $\kappa^* = \kappa_1 > \kappa_{min}$ .

*Proof.* First note that due to the assumptions,  $\psi^2$  is given by (28). Let us start by analyzing the (simpler) case  $-1/\gamma < -w^{(j)}$ , eliminating two of the subcases in (28). In view of the behavior of  $\psi^2(\kappa)$  as  $\kappa$  approaches infinity, we may assume that

$$\psi^2(\kappa) = (1 + \kappa) \left( \alpha - \frac{\beta^2 \kappa}{1 + \gamma \kappa} \right) = \frac{1 + \kappa}{1 + \gamma \kappa} [(\alpha\gamma - \beta^2)\kappa + \alpha], \quad (35)$$

and work with  $\kappa \in [-w^{(j)}, \infty)$ . If we discover that the infimum is attained “at”  $\infty$ , we will set  $\kappa^* = \infty$ . Consider the following cases:

1. If  $\beta = 0$  then  $\psi^2(\kappa) = (1 + \kappa)\alpha$ , which is nondecreasing, and hence  $\kappa^* = \kappa_{min}$ .
2. Assume that  $\beta \neq 0$  and notice that

$$(\psi^2)'(\kappa) = \alpha - \beta^2 \frac{\gamma\kappa^2 + 2\kappa + 1}{(1 + \gamma\kappa)^2} = \frac{\gamma(\alpha\gamma - \beta^2)\kappa^2 + 2(\alpha\gamma - \beta^2)\kappa + \alpha - \beta^2}{(1 + \gamma\kappa)^2}. \quad (36)$$

- (a) Let us first consider the degenerate case when the numerator in the expression above fails to be a quadratic. If  $\alpha\gamma = \beta^2$ , we see from (36) that  $\psi^2$  is increasing if  $\gamma < 1$  and hence we can choose  $\kappa^* = \kappa_{min}$ . If  $\gamma = 1$  then  $\psi^2$  is constant on  $[-w^{(j)}, \infty]$  and any choice of  $\kappa^*$  is optimal. Finally, if  $\gamma > 1$  then  $\kappa^* = \infty$ .

- (b) Assume that  $\alpha\gamma > \beta^2$ . The discriminant of the (convex) quadratic in the numerator of (36) is  $D = 4(\alpha\gamma - \beta^2)\beta^2(\gamma - 1)$ . This is nonpositive if  $\gamma \leq 1$ , in which case the derivative of  $\psi^2$  is nonnegative on  $(-1/\gamma, \infty) \supset [-w^{(j)}, \infty)$ . We can therefore choose  $\kappa^* = \kappa_{min}$ . Henceforth suppose  $\gamma > 1$  and let us write down the roots of the quadratic:

$$\kappa_{1,2} = \frac{-(\alpha\gamma - \beta^2) \pm |\beta|\sqrt{(\gamma - 1)(\alpha\gamma - \beta^2)}}{\gamma(\alpha\gamma - \beta^2)}.$$

Notice that

$$\kappa_2 = -\frac{1}{\gamma} - \frac{|\beta|\sqrt{\gamma - 1}}{\gamma\sqrt{\alpha\gamma - \beta^2}} < -\frac{1}{\gamma} < -\frac{1}{\gamma} + \frac{|\beta|\sqrt{\gamma - 1}}{\gamma\sqrt{\alpha\gamma - \beta^2}} = \kappa_1.$$

This implies that  $\psi^2$  is decreasing on  $(-1/\gamma, \kappa_1)$  and then increasing on  $(\kappa_1, \infty)$ . Since we consider only  $\kappa \geq -w^{(j)}$ , it is clear that  $\kappa^* = \max\{\kappa_{min}, \kappa_1\}$ .

It remains to analyze the situation with  $-1/\gamma = -w^{(j)}$ . In this case we proceed as above, except we have to take into account also the second and third expression in (28) defining  $\psi^2$ . If  $\beta = 0$  then  $\psi^2(\kappa) = (1 + \kappa)\alpha$  on  $[-w^{(j)}, \infty)$  and hence we conclude, as above, that  $\kappa^* = \kappa_{min}$ . Assume henceforth that  $\beta \neq 0$ . Now because  $\psi^2(\kappa) \rightarrow \infty = \psi^2(-w^{(j)})$  as  $\kappa \downarrow -w^{(j)}$ , we may proceed exactly as in the detailed analysis above, keeping in mind that  $\gamma > 1$ , which is a consequence of the assumption  $-1 < -w^{(j)} = -1/\gamma$ . In case 2a this leads to  $\kappa^* = \infty$ , while in case 2b we now know that  $-w^{(j)} = -1/\gamma < \kappa_1$  and hence  $\kappa^* = \kappa_1$ .  $\square$

## 6.2 Specialized line-search formula

Our initial motivation for the choice of  $j$  can be drawn from the multiplicative weight-update rule. At every iteration of Algorithm 1, the weights  $w^{(i)}$  are multiplied by the factor  $|\langle a_i, y \rangle|$  and then re-normalized. If this value is relatively large (or small) for particular  $i$ , the corresponding weight is being updated by a relatively large (or small) factor and is likely to have a substantial effect. It therefore makes sense to consider

$$j^+ := \arg \max_i |\langle a_i, y \rangle| \quad \text{and} \quad j^- := \arg \min_{w^{(i)} > 0} |\langle a_i, y \rangle|. \quad (37)$$

Notice that  $\varphi(y) = |\langle a_{j^+}, y \rangle|$  and that either  $a_j \in \partial\varphi(y)$  or  $-a_j \in \partial\varphi(y)$ , depending on whether  $\varphi(y) = \langle a_j, y \rangle$  or  $\varphi(y) = \langle -a_j, y \rangle$ .

If we assume that  $j$  is chosen to be either  $j^+$  or  $j^-$ , we can get a refined version of the optimal line-search formula. Let us first observe that  $\langle a_{j^-}, y \rangle^2 \leq \alpha \leq \langle a_{j^+}, y \rangle^2 = \varphi^2(y)$ , which is a simple consequence of the definitions of  $j^+$  and  $j^-$  and the frequently used identity  $\sum w^{(i)} \langle a_i, y \rangle^2 = \langle G(w)y, y \rangle = \langle d, y \rangle = \psi^2(w) = \alpha$ . Indeed, the above inequalities say that the weighted average of the numbers  $\langle a_i, y \rangle^2$  with positive weights  $w^{(i)}$  cannot be smaller than their minimum or bigger than their maximum. If there is equality in any of the two inequalities, then  $\langle a_i, y \rangle^2 = \alpha = \varphi^2(y)$

for all  $i$  for which  $w^{(i)} > 0$ , which is equivalent to the optimality condition (ii') of Theorem 20. So unless the current vector of weights  $w$  is optimal, we have

$$\langle a_{j-}, y \rangle^2 < \alpha < \langle a_{j+}, y \rangle^2 = \varphi^2(y). \quad (38)$$

Consider now the following cases:

1. Assume  $j = j^+$ . First notice that  $\alpha \leq \beta^2$ , with equality if and only if  $w$  is optimal. The Cauchy-Schwarz inequality  $\alpha\gamma \geq \beta^2$  then implies  $\gamma \geq 1$  and hence  $\gamma = 1$  implies optimality. Assume therefore that  $\gamma > 1$ , which excludes the first case in (33), and consider two subcases:

- (a) Case  $\alpha\gamma > \beta^2$ . By (33) we have  $\kappa^* = \max\{\kappa_{min}, \kappa_1\}$ . However, we can say a bit more. Noting that  $\alpha \leq \beta^2$  is equivalent to  $\kappa_1 \geq 0$ , we obtain  $\kappa^* = \kappa_1$ .
- (b) Case  $\alpha\gamma = \beta^2$ . Formula (33) implies  $\kappa^* = \infty$ . We claim that the next iterate (after taking the ‘infinite’ step) will be optimal. Indeed,  $G^+ := G(\kappa^*) = a_j a_j^*$  and if we let  $y^+$  satisfy  $G^+ y^+ = d$ , then

$$\sqrt{\alpha^+} := \|d\|_{G^+}^* = \langle d, y^+ \rangle^{1/2} = \langle a_j a_j^* y^+, y^+ \rangle^{1/2} = |\langle a_j, y^+ \rangle| = \frac{1}{\varphi^*}.$$

The last equality follows from  $\text{bdry } Q \ni \varphi^* d = \varphi^* \langle a_j, y^+ \rangle a_j$  because  $\{a_j, -a_j\} \subset \text{bdry } Q$  and hence it must be the case that  $|\varphi^* \langle a_j, y^+ \rangle| = 1$ .

2. Assume  $j = j^-$ . First note that  $\beta^2 = \langle a_{j-}, y \rangle^2 \leq \alpha$  with equality if and only if  $w$  is optimal.

If  $\gamma \leq 1$  then (33) implies  $\kappa^* = \kappa_{min}$ . If  $\gamma > 1$ , we get  $\beta^2 \leq \alpha < \alpha\gamma$  and consequently  $\kappa^* = \max\{\kappa_{min}, \kappa_1\}$ . Moreover, it is easy to show that  $\beta^2 \leq \alpha$  is equivalent to  $\kappa_1 \leq 0$ , which leads to the observation that  $\kappa^* \leq 0$ . If the current iterate is not optimal, then  $\beta^2 < \alpha$  and thus  $\kappa^* < 0$ .

We have arrived at the following conclusion:

**Theorem 30.** *Under the assumptions of Theorem 29 the following hold:*

1. If  $j = j^+$  then

$$\kappa^* = \begin{cases} \kappa_1 \geq 0 & \text{if } \gamma > 1 \text{ and } \alpha\gamma > \beta^2, \\ \infty & \text{if } \gamma > 1 \text{ and } \alpha\gamma = \beta^2. \end{cases} \quad (39)$$

Moreover, it is always the case that  $\alpha \leq \beta^2$ , with equality if and only if  $w$  is optimal. This happens, in particular, if  $\gamma = 1$ . The new iterate after the  $\kappa^* = \infty$  step is taken is optimal.

2. If  $j = j^-$  then

$$\kappa^* = \begin{cases} \kappa_{min} & \text{if } \gamma \leq 1, \\ \max\{\kappa_{min}, \kappa_1\} \leq 0 & \text{if } \gamma > 1. \end{cases} \quad (40)$$

Moreover, it is always the case that  $\beta^2 \leq \alpha$ , with equality if and only if  $w$  is optimal.

## 7 An algorithm with “increase” steps only

In this subsection we design and analyze an algorithm which at every iteration uses the choice  $j = j^+ = \arg \max_i |\langle a_i, y \rangle|$ , where  $y$  is some vector satisfying  $Gy = d$ , and updates  $G$  to  $G(\kappa) = (G + \kappa a_j a_j^*) / (1 + \kappa)$ , using the optimal step size  $\kappa^*$  described by Corollary 30. Since this particular choice of  $j$  always leads to nonnegative value of the optimal step size parameter, strictly positive if  $w$  is not optimal, we see from the definition of  $w(\kappa)$  (22) that the weight  $w^{(j)}$  will increase while all other weights decrease uniformly to account for this. This explains the choice of the terminology “increase” step.

Since the initial iterate  $w_0$  used in Algorithm 2 has all components positive (all are equal to  $\frac{1}{m}$ ), all weights stay positive throughout the algorithm. In other words, the method proceeds through the interior of the feasible region. One important consequence of this is that the iterate matrices  $G$  never lose rank and hence stay positive definite throughout the algorithm. This implies that a system of the form  $Gy = d$  will always have a unique solution, which is the first step towards an implementable code. Of course, numerical instabilities might occur in situations when certain weights get close to zero and, as a result,  $G$  becomes nearly rank-deficient (see (10)). In this work we do not present any strategies for dealing with this linear algebra issue and instead focus on the optimization-theoretic results. Let us remark that it is unlikely that there will be problems with solving  $Gy = d$  for a general position of the vector  $d$ .

The analysis of Algorithm 2 is based on a result which says that a certain approximation of the optimal step size gives a sufficient decrease in the value of  $\psi^2$  (see Lemma 31 below). Let us first describe a motivational heuristic, leading us to the discovery of a suitable approximately optimal step size.

### 7.1 A step-size heuristic

Let

$$\hat{\delta} := \frac{|\beta|}{\sqrt{\alpha}} - 1, \quad (41)$$

which can be also written as  $\beta^2 = \alpha(1 + \hat{\delta})^2$ , and assume  $\hat{\delta} > 0$ . By Theorem 21, the current iterate  $w$  is  $\hat{\delta}$ -optimal for (P3). To see this, we just need to translate our simplified notation (using  $\alpha, \beta$  and  $\gamma$ ) to the symbols used in that theorem:  $\varphi(y) = |\beta|$  and  $\psi(w) = \sqrt{\alpha}$ . We now see that the definition of  $\hat{\delta}$  implies  $\varphi(y) = (1 + \hat{\delta})\psi(w)$ , which is precisely the universal  $\hat{\delta}$ -approximate optimality condition (20), implying  $\psi(w) \leq (1 + \hat{\delta})\psi^*$ .

Assume, just for the sake of the motivational heuristic to follow, that  $\alpha\gamma > \beta^2$ . This excludes the “next-iterate is optimal” case from the description of the optimal step size of Corollary 30 and implies  $\kappa^* = \kappa_1$ . We claim that in the situation when  $\gamma$  is large,  $\kappa^*$  is reasonably well approximated by  $\hat{\delta}/\gamma$ . Indeed, the ratio  $\kappa_1/(\hat{\delta}/\gamma)$  converges to 1 (from above) as  $\gamma$  approaches

infinity:

$$\frac{\kappa_1}{\frac{\hat{\delta}}{\gamma}} = \frac{-\frac{1}{\gamma} + \frac{|\beta|\sqrt{\gamma-1}}{\gamma\sqrt{\alpha\gamma-\beta^2}}}{\frac{|\beta|-1}{\sqrt{\alpha}}}{\frac{|\beta|\sqrt{\frac{\gamma-1}{\alpha\gamma-\beta^2}} - 1}{|\beta|\sqrt{\frac{1}{\alpha}} - 1}} \downarrow 1 \quad \text{as} \quad \gamma \rightarrow \infty.$$

Note that the convergence is “from above” because  $(\gamma-1)/(\alpha\gamma-\beta^2) > 1/\alpha$ , which follows from  $\beta^2 > \alpha$ .

## 7.2 Sufficient decrease

Being motivated by the optimal step-size approximation discussed above, we now show that if the condition for  $\delta$ -approximate optimality of Theorem 21 is not met, then by taking the step  $\delta/\gamma$ , we can reduce the (square of the) objective value by at least  $\alpha\delta^2/\gamma$ . This will play a central role in the analysis of our algorithm. Note that by taking a finite step ( $\kappa \neq \infty$ ), the function  $\psi^2$  decreases by

$$\varepsilon(\kappa) := \psi^2(0) - \psi^2(\kappa) = \frac{\beta^2\kappa(1+\kappa)}{1+\gamma\kappa} - \alpha\kappa. \quad (42)$$

If  $\kappa = \infty$  is used, this formula should be understood in the limit sense – see (29).

**Lemma 31** (Sufficient decrease). *If  $|\beta| \geq (1+\delta)\sqrt{\alpha}$  for some  $\delta \geq 0$  and we let  $\kappa := \frac{\delta}{\gamma}$ , then*

$$\varepsilon(\kappa^*) \geq \varepsilon(\kappa) \geq \alpha\frac{\delta^2}{\gamma}.$$

*Proof.*

$$\begin{aligned} \varepsilon(\kappa^*) \geq \varepsilon(\kappa) &= \frac{\beta^2\kappa(1+\kappa)}{1+\gamma\kappa} - \alpha\kappa \geq \frac{(1+\delta)^2\alpha\frac{\delta}{\gamma}(1+\frac{\delta}{\gamma}) - (1+\delta)\alpha\frac{\delta}{\gamma}}{1+\delta} \\ &= \frac{\alpha\delta}{\gamma}[(1+\delta)(1+\frac{\delta}{\gamma}) - 1] \\ &\geq \frac{\alpha\delta^2}{\gamma}. \end{aligned}$$

The last inequality follows from the estimate  $1 + \frac{\delta}{\gamma} \geq 1$ . □

**Remark 32.** *Observe that the condition of the above lemma is satisfied with equality for  $\delta = \hat{\delta}$  defined in (41).*

### 7.3 Even better decrease

It turns out that if we take into account also some other choices of  $j$ , we can possibly achieve an even bigger decrease in  $\psi^2$  than that guaranteed by the above lemma. Let  $\beta^i := \langle a_i, y \rangle$  and  $\gamma^i := \langle a_i, G^{-1}a_i \rangle$  (for all  $i$ ), so that  $\beta^j = \beta$  and  $\gamma^j = \gamma$ . Also let  $\delta^i := (|\beta^i|/\sqrt{\alpha}) - 1$  and  $\mathcal{I}$  be the set of those indices  $i$  for which  $\delta^i \geq 0$ . Note that  $j = j^+ \in \mathcal{I}$ . Observe that the argument of the above lemma can be repeated to show that

$$\varepsilon(\kappa^*) \geq \varepsilon(\delta^i/\gamma^i) \geq \alpha \frac{(\delta^i)^2}{\gamma^i} = \frac{(|\beta^i| - \sqrt{\alpha})^2}{\gamma^i}, \quad \forall i \in \mathcal{I}. \quad (43)$$

Note that this is the lower bound on the change in  $\psi^2$  when using  $a_i$  instead of  $a_j$  and the corresponding approximately optimal step size. While the specific choice  $i = j$  guarantees sufficient decrease, we might be able to do better by *optimizing* over the set  $\mathcal{I}$  rather than by picking the *feasible* solution  $i = j \in \mathcal{I}$ . This leads us to defining

$$i^* := \arg \max_{i \in \mathcal{I}} \frac{(|\beta^i| - \sqrt{\alpha})^2}{\gamma^i}.$$

Certainly, the *decrease guaranteed* by  $i^*$  is at least as good as the decrease guaranteed by  $j$ . This does not mean, however, that the *actual decrease*, by taking the optimal step, will be bigger. This has to be taken into account when implementing this strategy in an algorithm. If we decide to use this improvement, we have to deal with the issue of the actual computation of the values  $\gamma^i$ , which are needed to find  $i^*$ . By doing the computation from scratch at every iteration, we will need  $O(mn^2)$  arithmetic operations:  $O(n^2)$  for solving each of the at most  $m$  equations  $Gx = a_i$  ( $i \in \mathcal{I}$ ), assuming we maintain the Cholesky factorization of  $G$  from iteration to iteration. Alternatively, we can solve for  $G_0^{-1}a_i$  for all  $i$  at the beginning of the algorithm, which takes  $O(mn^2)$  operations if we assume the availability of the Cholesky factors of  $G_0$ , and subsequently update the solutions as we modify the (Cholesky factorization of the) matrix. The work per iteration will drop to  $O(mn)$ , which is of the same order as the work needed to calculate  $j$ .

### 7.4 A crucial assumption

**Assumption 33.** *The values  $\gamma_k$  generated by Algorithm 2 are bounded above by some constant  $\Gamma$ .*

As we shall see, the parameter  $\Gamma$  appears in the complexity bounds. It would therefore be good to be able to estimate its size. We will deal with this issue in Subsection 9.

### 7.5 Quick and dirty analysis

Let us first offer a rough analysis of Algorithm 2, leading to the performance guarantee of  $O(\Gamma\delta^{-2} \ln m)$  iterations of a first-order method, followed by a more refined analysis with the

---

**Algorithm 2 (Inc)** Solving (P3) using increase steps only.

---

- 1: **Input:**  $a_1, \dots, a_m \in \mathbf{E}^*$ ,  $d \in \mathbf{E}^*$ ,  $\delta > 0$ ;
  - 2: **Initialize:**  $k = 0$ ,  $w_0 = (1/m)e_m$ ,  $G_0 = \frac{1}{m} \sum_i a_i a_i^*$ ,  $y_0 = G_0^{-1} d$ ;
  - 3: **Iterate:**
  - 4:    $\alpha_k = \langle d, y_k \rangle$ ,  $j = \arg \max_i |\langle a_i, y_k \rangle|$ ,  $g_k = a_j$ ;
  - 5:    $\beta_k = \langle g_k, y_k \rangle$ ,  $\gamma_k = \langle g_k, G_k^{-1} g_k \rangle$ ;
  - 6:    $\delta_k = \frac{|\beta_k|}{\sqrt{\alpha_k}} - 1$ ;
  - 7:   **if**  $\delta_k \leq \delta$
  - 8:     **terminate**;
  - 9:   **else**
  - 10:    **if**  $\alpha_k \gamma_k = \beta_k^2$
  - 11:      $\kappa^* = \infty$ ,  $G_{k+1} = g_j g_j^*$ ,  $w_{k+1} = \frac{w_k + \kappa^* e_j}{1 + \kappa^*} = \underbrace{(0, \dots, 0)}_{1 \dots j-1}, \underbrace{1}_j, \underbrace{(0, \dots, 0)}_{j+1 \dots m}$ ;
  - 12:      $k \leftarrow k + 1$ , **terminate**; ( $w_k$  is optimal)
  - 13:    **else**
  - 14:      $\kappa^* = -\frac{1}{\gamma_k} + \frac{|\beta_k| \sqrt{\gamma_k - 1}}{\gamma_k \sqrt{\alpha_k \gamma_k - \beta_k^2}}$ ,  $G_{k+1} = \frac{G_k + \kappa^* g_k g_k^*}{1 + \kappa^*}$ ,  $w_{k+1} = \frac{w_k + \kappa^* e_j}{1 + \kappa^*}$ ;
  - 15:      $y_{k+1} = (1 + \kappa^*) \left( y_k - \frac{\beta_k \kappa^* G_k^{-1} g_k}{1 + \gamma_k \kappa^*} \right)$ ;
  - 16:      $k \leftarrow k + 1$ ;
  - 17:    **end if**
  - 18: **Output:**  $w_k$  satisfying  $\|d\|_{G_k}^* = \sqrt{\alpha_k} \leq (1 + \delta) \psi^*$  and  $G_k = G(w_k)$
-

guarantee  $O(\Gamma(\ln \Gamma + \ln \ln m + \delta^{-1}))$ . For the quick result observe that because

$$\varphi^* d \in Q \subseteq \sqrt{m} \mathcal{B}(G_0), \quad (44)$$

we have  $\sqrt{\alpha_0} = \|d\|_{G_0}^* \leq \sqrt{m}/\varphi^*$ . Now assume that Algorithm 2 produces  $K + 1$  iterates with  $K + 1 \geq \lceil \Gamma \delta^{-2} \ln m \rceil$ . The termination criterion of line 7 then implies that  $\delta_k > \delta$  for  $k = 0, 1, \dots, K$ . Since  $|\beta_k| = (1 + \delta_k)\sqrt{\alpha_k}$  for all  $k \leq K$ , Lemma 31 and Assumption 33 imply

$$\alpha_k - \alpha_{k+1} \geq \alpha_k \frac{\delta_k^2}{\gamma_k} > \alpha_k \frac{\delta^2}{\Gamma}.$$

Repeated use of this inequality gives  $\alpha_{K+1} < \alpha_0(1 - \delta^2/\Gamma)^{K+1}$  and hence

$$\begin{aligned} \psi^2(w_{K+1}) &= \alpha_{K+1} < \alpha_0(1 - \delta^2/\Gamma)^{K+1} \\ &\leq \alpha_0 e^{-(K+1)\delta^2/\Gamma} \leq \alpha_0 e^{-\ln m} \leq \frac{m}{(\varphi^*)^2} \frac{1}{m} = (\psi^*)^2, \end{aligned}$$

which contradicts the fact that  $\psi^*$  is the optimal value of problem (P3).

## 7.6 Refined analysis

The following is the central result of this paper:

**Theorem 34.** *Under Assumption 33, Algorithm 2 produces a  $\delta$ -approximate solution of (P3) (and hence by Theorem 21 of (P1), (D1), (D'1), (P2), (D2) and (D3)) in at most*

$$2\Gamma \left( \ln \Gamma + \ln \ln m + \frac{8}{\delta} \right)$$

*iterations.*

*Proof.* Let  $L_k := \ln \sqrt{\alpha_k}$  and  $L^* = \ln \psi^*$  and notice that  $\sqrt{\alpha_k} \leq (1 + \delta_k)\psi^*$ . By taking logarithms,

$$\varepsilon'_k := L_k - L^* \leq \ln(1 + \delta_k). \quad (45)$$

Also,  $\beta_0^2 = \max_i \langle a_i, y_0 \rangle^2 \leq \sum_i \langle a_i, y_0 \rangle^2 = m \langle G_0 y_0, y_0 \rangle = m \alpha_0 = m \frac{\beta_0^2}{(1 + \delta_0)^2}$ , whence

$$\delta_0 \leq \sqrt{m} - 1 \quad \text{and} \quad \varepsilon'_0 \leq \ln(1 + \delta_0) \leq \frac{1}{2} \ln m. \quad (46)$$

By Lemma 31,  $\varepsilon(\kappa^*) = \alpha_k - \alpha_{k+1} \geq \alpha_k \delta_k^2 / \gamma_k \geq \alpha_k \delta_k^2 / \Gamma$  and therefore

$$\alpha_{k+1} \leq \alpha_k (1 - \delta_k^2 / \Gamma). \quad (47)$$

By taking logarithms in (47) and using (45),

$$L_k - L_{k+1} \geq -\frac{1}{2} \ln(1 - \delta_k^2 / \Gamma) \geq \frac{1}{2} \delta_k^2 / \Gamma \geq \frac{1}{2\Gamma} \ln(1 + \delta_k^2) \geq \frac{1}{2\Gamma} \ln(1 + \delta_k), \quad (48)$$

with the last inequality true whenever  $\delta_k \geq 1$ . Combining (45) and (48) yields

$$\varepsilon'_{k+1} \leq \varepsilon'_k \left(1 - \frac{1}{2\Gamma}\right),$$

for all  $k$  with  $\delta_k \geq 1$ . We will now bound the number of iterations for which  $\delta_k \geq 1$ . The last inequality together with (46) gives

$$\varepsilon'_k \leq \varepsilon'_0 \left(1 - \frac{1}{2\Gamma}\right)^k \leq \frac{1}{2} \ln m \exp\left(-\frac{k}{2\Gamma}\right). \quad (49)$$

Due to (49) and  $\varepsilon'_k \geq \varepsilon'_k - \varepsilon'_{k+1} = L_k - L_{k+1} \geq \frac{1}{2}\delta_k^2/\Gamma \geq \frac{1}{2}\Gamma^{-1}$ , the largest  $k$  for which  $\delta_k \geq 1$  must satisfy  $\Gamma^{-1} \leq \ln m \exp\left(-\frac{k}{2\Gamma}\right)$ , leading to the bound

$$k \leq 2\Gamma(\ln \Gamma + \ln \ln m). \quad (50)$$

So one can obtain a solution within the factor of 2 of the optimum in  $O(\Gamma(\ln \Gamma + \ln \ln m))$  iterations of Algorithm 2.

Following the “halving” argument of Khachiyan [15], we can bound the number of additional iterations needed to obtain the desired  $\delta$ -approximate solution. Suppose  $\delta_k \leq 1$ , and let  $h(\delta_k)$  be the smallest integer  $h$  such that  $\delta_{k+h} \leq \delta_k/2$ . Whenever  $\delta_{k+h} \geq \delta_k/2$ , we also have

$$\varepsilon'_{k+h} - \varepsilon'_{k+h+1} \geq \frac{1}{2}\delta_{k+h}^2/\Gamma \geq \frac{1}{8}\delta_k^2/\Gamma,$$

which says that the gap in (45) must at every such iteration decrease by at least  $\frac{1}{8}\delta_k^2/\Gamma$ . However, the original gap is of size at most  $\varepsilon'_k \leq \ln(1 + \delta_k) \leq \delta_k$ , and hence the number of iterations needed for halving  $\delta_k$  is bounded above by

$$h(\delta_k) \leq \frac{\delta_k}{\frac{1}{8}\delta_k^2/\Gamma} = \frac{8\Gamma}{\delta_k}.$$

In order to get below  $\delta$ , we need to “halve”  $l$ -times where  $l$  is obtained from  $\delta_k/2^l \leq \delta$ , that is  $l = \lceil \log_2 \delta_k/\delta \rceil$ , where  $k$  is the first iteration for which  $\delta_k \leq 1$ . The total number of additional iterations required to achieve the desired  $\delta$ -approximate solution is at most

$$\sum_{i=0}^{l-1} h(\delta_k/2^i) \leq 8\Gamma \sum_{i=0}^{l-1} \frac{1}{\delta_k/2^i} = \frac{8\Gamma}{\delta_k} 2^{\lceil \log_2 \delta_k/\delta \rceil} \leq \frac{16\Gamma}{\delta}.$$

□

## 8 An algorithm with “increase”, “decrease” and “drop” steps

In the previous subsection we have analyzed an algorithm which at every iteration works with  $j = j^+$  (an “increase” step). A consequence of this choice is that the optimal step-size parameter  $\kappa^*$  is always nonnegative, implying that  $w^{(j)}$  is being increased while all other weights are decreased uniformly (and hence at a slower rate than the rate of increase of  $w^{(j)}$ ) in due compensation.

Starting from  $w_0 = (\frac{1}{m}, \dots, \frac{1}{m})$ , Algorithm 2 keeps all weights positive until termination. In an optimal solution  $w$ , however, the weights can be positive only for points  $a_i$  lying on a face (say  $\mathcal{F}$ ) of  $Q$  containing the point  $\varphi^*d$  — the intersection of  $Q$  and the half-line emanating from the origin in the direction  $d$  (see Figure 2). Note that in the case when  $m \gg n$ , it is to be expected that many more points will have zero weights rather than positive weights, at optimality. It therefore seems intuitive that if the incorporation of “decrease” and/or “drop” steps could speed up the algorithm considerably.

In this subsection we propose and analyze an algorithm in which we allow also for “decrease” and “drop” iterations — steps which decrease  $w^{(j)}$ , respectively drop it to zero ( $\kappa = -w^{(j)}$ ). The idea is as follows. At every iteration we consider both  $j = j^+$  and  $j = j^-$ . We make the latter choice if the predicted decrease is better (this corresponds to  $\delta^- \geq \delta^+$  in Algorithm 3), except when this leads to a drop step reducing the rank of  $G$  (this happens when  $-\frac{1}{\gamma} = -w^{(j)} = \kappa$ ). Otherwise we choose  $j = j^+$ .

There are several reasonable alternative rules for deciding among  $j^+$  and  $j^-$ . For example, we could base our decision on comparing the *actual decrease* as opposed to the *decrease predicted* by  $\delta^+$  and  $\delta^-$ . We could also forbid taking drop steps altogether, allowing only for decrease steps, etc.

Let us start with a twin result to Lemma 31 which essentially says that if we choose  $j = j^-$  and it happens that  $\kappa^*$  is not a drop step, then by taking this step we are guaranteed sufficient decrease in the (square of the) objective function:

**Lemma 35.** *Assume  $j = j^-$ .*

(i) *If  $|\beta| \leq (1 - \delta)\sqrt{\alpha}$  for some  $0 \leq \delta < 1$  and  $\kappa := -\frac{\delta}{\gamma} \geq -w^{(j)}$ , then*

$$\varepsilon(\kappa^*) \geq \varepsilon(\kappa) \geq \alpha \frac{\delta^2}{\gamma}.$$

(ii) *If  $\kappa^* = \kappa_1$  and  $\delta := 1 - \frac{|\beta|}{\sqrt{\alpha}}$ , then  $\kappa := -\frac{\delta}{\gamma} \geq -w^{(j)}$ .*

*Proof.* For part (i) notice that the assumption  $\kappa \geq -w^{(j)} = \kappa_{min}$  ensures feasibility of the line-search parameter  $\kappa$ . Also observe that  $1 - \frac{\delta}{\gamma} \geq 1 - w^{(j)} > 0$  by Assumption 2 in Subsection 5.2. We now proceed as in Lemma 31:

$$\begin{aligned} \varepsilon(\kappa^*) \geq \varepsilon(\kappa) &= \frac{\beta^2 \kappa (1 + \kappa)}{1 + \gamma \kappa} - \alpha \kappa \\ &\geq \frac{-(1 - \delta)^2 \alpha \frac{\delta}{\gamma} (1 - \frac{\delta}{\gamma}) + (1 - \delta) \alpha \frac{\delta}{\gamma}}{1 - \delta} \\ &= \frac{\alpha \delta}{\gamma} [1 - (1 - \delta)(1 - \frac{\delta}{\gamma})] \\ &\geq \frac{\alpha \delta^2}{\gamma}. \end{aligned}$$

---

**Algorithm 3 (IncDec)** Solving (P3) using both increase and decrease steps.

---

- 1: **Input:**  $a_1, \dots, a_m \in \mathbf{E}^*$ ,  $d \in \mathbf{E}^*$ ,  $\delta > 0$ ;
  - 2: **Initialize:**  $k = 0$ ,  $w_0 = (1/m)e_m$ ,  $G_0 = \frac{1}{m} \sum_i a_i a_i^*$ ,  $y_0 = G_0^{-1} d$ ;
  - 3: **Iterate:**
  - 4:    $\alpha_k = \langle d, y_k \rangle$ ;
  - 5:    $j^- = \arg \min_i \{ |\langle a_i, y_k \rangle| : w_k^{(i)} > 0 \}$ ,  $\beta^- = \langle a_{j^-}, y_k \rangle$ ,  $\delta^- = 1 - \frac{|\beta^-|}{\sqrt{\alpha_k}}$ ;
  - 6:    $j^+ = \arg \max_i |\langle a_i, y_k \rangle|$ ,  $\beta^+ = \langle a_{j^+}, y_k \rangle$ ,  $\delta^+ = \frac{|\beta^+|}{\sqrt{\alpha_k}} - 1$ ;
  - 7:   **if**  $\delta^+ \leq \delta$  **then terminate; end if**
  - 8:   **if**  $\delta^+ < \delta^-$
  - 9:      $j = j^-$ ,  $g_k = a_j$ ,  $\beta_k = \beta^-$ ,  $\gamma_k = \langle g_k, G_k^{-1} g_k \rangle$ ,  $\delta_k = \delta^-$ ;
  - 10:    **if**  $\gamma_k > 1$    **then**    $\kappa = \max \left\{ -\frac{1}{\gamma_k} + \frac{|\beta_k| \sqrt{\gamma_k - 1}}{\gamma_k \sqrt{\alpha_k \gamma_k - \beta_k^2}}, -w_k^{(j)} \right\}$ ;
  - 11:     **else**    $\kappa = -w_k^{(j)}$ ;   **end if**
  - 12:    **if**  $\kappa = -w_k^{(j)} = -\frac{1}{\gamma_k}$  **then jump to 14, end if**
  - 13:    **else**
  - 14:      $j = j^+$ ,  $g_k = a_j$ ,  $\beta_k = \beta^+$ ,  $\gamma_k = \langle g_k, G_k^{-1} g_k \rangle$ ,  $\delta_k = \delta^+$ ;
  - 15:     **if**  $\alpha_k \gamma_k > \beta_k^2$    **then**    $\kappa = -\frac{1}{\gamma_k} + \frac{|\beta_k| \sqrt{\gamma_k - 1}}{\gamma_k \sqrt{\alpha_k \gamma_k - \beta_k^2}}$ ;
  - 16:     **else**    $\kappa = \infty$ ;   (the next iterate is optimal) **end if**
  - 17:    **end if**
  - 18:     $w_{k+1} = \frac{w_k + \kappa e_j}{1 + \kappa}$ ,  $G_{k+1} = \frac{G_k + \kappa g_k g_k^*}{1 + \kappa}$ ,  $y_{k+1} = (1 + \kappa) \left( y_k - \frac{\beta_k \kappa G_k^{-1} g_k}{1 + \gamma_k \kappa} \right)$ ;
  - 19:     $k \leftarrow k + 1$ ;
  - 20: **Output:**  $w_k$  satisfying  $\|d\|_{G_k}^* = \sqrt{\alpha_k} \leq (1 + \delta) \psi^*$  and  $G_k = G(w_k)$
-

The last inequality follows from the estimate  $0 < 1 - \frac{\delta}{\gamma} \leq 1$ . Let us now prove (ii). Because  $\kappa^* = \kappa_1$  is feasible for the line-search problem, we must have  $\kappa_1 \geq -w^{(j)}$ . However, using the inequality  $\beta^2 \leq \alpha$  it can be argued by simple algebra that  $-\frac{\delta}{\gamma} \geq \kappa_1$  (see (34) for the definition of  $\kappa_1$ ).  $\square$

**Theorem 36.** *Under Assumption 33, Algorithm 3 produces a  $\delta$ -approximate solution of (P3) (and hence by Theorem 21 of (P1), (D1), (D'1), (P2), (D2) and (D3)) in at most*

$$m + 4\Gamma \left( \ln \Gamma + \ln \ln m + \frac{8}{\delta} \right)$$

*iterations.*

*Proof.* Due to Lemma 35, the argument is identical to the proof of Theorem 34. The difference is that we need to bound the number of drop iterations because these do not guarantee any positive decrease (but do not increase the objective either). Note that either the current point  $a_j$  is dropped for the first time (there are a maximum of  $m$  such occurrences), or it has been dropped before, in which case we can pair it up with the previous iteration that increased the weight  $w^{(j)}$  from zero to a positive value. This algorithm therefore needs at most  $m$  plus twice the number of iterations guaranteed by Theorem 34.  $\square$

**Remark 37.** *The  $\ln \ln m$  factor in the complexity estimates of Algorithms 2 and 3 can be replaced by  $\ln \ln n$  if we pre-compute a rounding of  $Q$  with  $\frac{1}{\alpha} = O(\sqrt{n})$  and use the corresponding matrix as  $G_0$ . This can be done in  $O(n^2 m \log m)$  arithmetic operations (see [21]).*

## 9 Bounding the unknown constant

The performance guarantees of Algorithms 2 and 3 depend on the assumption that the squared norms of the points  $a_j$  encountered throughout the iterations are bounded from above by some constant  $\Gamma$ . It is therefore highly desirable to invest some time into exploring our options of theoretical and/or practical justification of this assumption.

How large can  $\Gamma$  be? Notice that we know from Corollary 26 that for any  $j$  with positive weight  $w^{(j)}$ , the value

$$\gamma_j := (\|a_j\|_{G(w)}^*)^2$$

can be bounded from above by a function of  $w^{(j)}$ :

$$\gamma_j \leq \frac{1}{w^{(j)}}. \tag{51}$$

If  $w^{(j)} = 0$ , as is the case when we perform an “add” step in Algorithm 3, we do not have an upper bound on  $\gamma_j$ . If we maintain all weights positive, as in Algorithm 2, then  $\Gamma$  can certainly be bounded by the reciprocal of the smallest weight  $w^{(j)}$  encountered throughout the algorithm. This leads to the idea of modifying our methods so as to keep all weights above a certain positive constant.

## 9.1 Bounding the weights away from zero: theoretical implications

Motivated by the above discussion, let us explicitly require that all weights be bounded away from zero by  $\frac{\varepsilon}{m}$ , with  $\varepsilon \in [0, 1]$  being a small constant *independent of the dimensions* of the problem. Note that setting  $\varepsilon = 1$  implies that all weights are equal to  $\frac{1}{m}$ .

It seems to be intuitively sound to expect that if we restrict the set of feasible points of problem (P3) by requiring  $w^{(i)} \geq \frac{\varepsilon}{m}$  for all  $i$ , the optimal value of the modified problem, which we will call  $(P3_\varepsilon)$ , should be close to the optimal value of (P3). Also, as  $\varepsilon$  gets smaller, the optimal value of  $(P3_\varepsilon)$  should approach that of (P3). We will formalize these ideas in the remainder of this subsection. Let

$$\Delta_m^\varepsilon := \{w \in \Delta_m : w^{(i)} \geq \frac{\varepsilon}{m}, i = 1, 2, \dots, m\}$$

and consider the following problem

$$(P3_\varepsilon) \quad \boxed{\psi_\varepsilon^* := \min_w \{\psi(w) : w \in \Delta_m^\varepsilon\}.}$$

We claim that the value  $\psi_\varepsilon^*$  is close to  $\psi^*$  for small  $\varepsilon$ :

**Theorem 38.** *For the optimal values  $\psi^*$  and  $\psi_\varepsilon^*$  of (P3) and  $(P3_\varepsilon)$ , respectively, we have*

$$\psi_\varepsilon^* \leq \frac{1}{(1 - \frac{m-1}{m}\varepsilon)^{1/2}} \psi^*. \quad (52)$$

To prove this we will need an auxiliary result.

**Lemma 39.** *For any  $x \in \mathbf{E}$ ,*

$$\max_{w \in \Delta_m^\varepsilon} \|x\|_{G(w)} \geq (1 - \frac{m-1}{m}\varepsilon)^{1/2} \varphi(x).$$

*Proof.* Assume  $\varphi(x) = |\langle a_j, x \rangle|$  and let  $\bar{w}$  be a vector of weights with  $\bar{w}^{(j)} = 1 - \frac{m-1}{m}\varepsilon$  and  $\bar{w}^{(i)} = \frac{1}{m}\varepsilon$  for all other  $i$ . Then

$$\max_{w \in \Delta_m^\varepsilon} \|x\|_{G(w)} \geq \|x\|_{G(\bar{w})} = \left( \sum \bar{w}^{(i)} \langle a_i, x \rangle^2 \right)^{1/2} \geq (\bar{w}^{(j)})^{1/2} |\langle a_j, x \rangle|.$$

□

*Proof.* (theorem)

$$\begin{aligned} \frac{1}{\psi_\varepsilon^*} &= \left[ \min_{w \in \Delta_m^\varepsilon} \|d\|_{G(w)}^* \right]^{-1} = \max_{w \in \Delta_m^\varepsilon} 1/\|d\|_{G(w)}^* \\ &= \max_{w \in \Delta_m^\varepsilon} \min_{\langle d, x \rangle=1} \|x\|_{G(w)} \\ &= \min_{\langle d, x \rangle=1} \max_{w \in \Delta_m^\varepsilon} \|x\|_{G(w)} \\ &\geq \min_{\langle d, x \rangle=1} (1 - \frac{m-1}{m}\varepsilon)^{1/2} \varphi(x) \\ &= (1 - \frac{m-1}{m}\varepsilon)^{1/2} \varphi^* = (1 - \frac{m-1}{m}\varepsilon)^{1/2} \frac{1}{\psi^*}. \end{aligned}$$

The exchange of the maximum and minimum can be justified by using Hartung's minimax theorem [14].  $\square$

**Remark 40.** For  $\varepsilon = 1$ , inequality (52) states that  $\psi(w_0) \leq \sqrt{m}\psi^*$ , where  $w_0$  is the vector of all weights equal to  $\frac{1}{m}$ . This we have already seen before as a consequence of the rounding property (44) of  $G_0 = G(w_0)$ .

**Corollary 41.** If  $\varepsilon \leq \frac{1}{2\tau}(\sqrt{\tau(\tau+4)} + \tau - 2)$  for some positive parameter  $\tau$  (necessarily,  $\tau \geq \frac{1}{2}$ ), then  $\psi_\varepsilon^* \leq (1 + \tau\varepsilon)\psi^*$ . In particular, if  $\varepsilon \leq \frac{1}{2}(\sqrt{5} - 1)$ , then  $\psi_\varepsilon^* \leq (1 + \varepsilon)\psi^*$ .

*Proof.* The condition on  $\varepsilon$  is equivalent to the last inequality in  $(1 - \frac{m-1}{m}\varepsilon)^{-1/2} \leq (1 - \varepsilon)^{-1/2} \leq (1 + \tau\varepsilon)$ .  $\square$

## 9.2 Bounding the weights away from zero: algorithmic implications

It is not trivial to see how one would go about modifying our algorithms to efficiently solve  $(P3_\varepsilon)$ . The requirement of keeping the weights above some positive threshold value  $\frac{\varepsilon}{m}$  does not seem to be cheap to maintain. Let us briefly explain why.

One possible approach to solving  $(P3_\varepsilon)$  using our methodology would involve dividing the operator  $G$  into two parts, keeping one fixed, ensuring that the weights are kept above  $\frac{\varepsilon}{m}$ . The other is a variable part, consisting of the remaining portion of the total weight. That is, we write

$$G = \sum_{i=1}^m \frac{\varepsilon}{m} a_i a_i^* + \sum_{i=1}^m \bar{w}^{(i)} a_i a_i^* = G_\varepsilon + G(\bar{w}),$$

where  $\sum_i \bar{w}^{(i)} = 1 - \varepsilon$ ,  $\bar{w}^{(i)} \geq 0$ ; that is,  $\bar{w} \in (1 - \varepsilon)\Delta_m$ . One would now update only the variable part, similarly as in the previous analysis:

$$G(\kappa) = G_\varepsilon + \frac{G(\bar{w}) + \kappa a_j a_j^*}{1 - \varepsilon + \kappa}. \quad (53)$$

Notice that we no longer have  $1 + \kappa$  in the denominator, and this would need to be accounted for by reworking the relevant analysis. The main problem with this approach, however, is that (53) no longer constitutes a simple enough update of the operator  $G$ . It is certainly not a rank-one-and-scaling update, as before. This means that it could be hard to be able to use the information from the previous iteration (for example, the Cholesky factor of  $G$  and the solution  $y$  of  $Gy = d$ ) to solve the new system  $G(\kappa)y = d$ . If we need to solve this from scratch, it requires  $O(n^3)$  arithmetic operations (assuming  $G(\kappa)$  is assembled from  $G_\varepsilon$  and  $G(\bar{w})$  via (53), which takes only  $O(n^2)$  arithmetic operations), which is worse than the previous  $O(n^2)$  work. However, the per-iteration arithmetical complexity of Algorithms 2 and 3 is  $O(mn)$ , which will dominate the work above in the case when  $m \geq n^2$ . The critical saving would then come from the fact that we do not have to form the new matrix from scratch, which would otherwise require  $O(mn^2)$  arithmetic operations.

While in this paper we do not show any details of a direct algorithm of this type for solving  $(P3_\varepsilon)$ , we believe that the ideas we have just described could be turned into a provably working algorithm, albeit one with a considerably higher computational effort per iteration.

### 9.3 The average of the gammas

As a possible alternative to the conservative strategy of keeping *all* weights above a certain positive threshold value throughout the algorithm, let us briefly discuss if it is possible to instead select a *particular*  $j$  so that  $\gamma_j$  is of a reasonable size. Let us start with the following simple observation:

**Lemma 42.**

$$\sum_{w^{(i)} > 0} w^{(i)} \gamma_i = \text{rank } G(w).$$

*Proof.* Assume first  $G := G(w)$  is invertible. Then

$$\begin{aligned} \sum_{w^{(i)} > 0} w^{(i)} \gamma_i &= \sum_i w^{(i)} \langle a_i, G(w)^{-1} a_i \rangle = \sum_i w^{(i)} \text{trace}[\langle a_i, G(w)^{-1} a_i \rangle] \\ &= \sum_i w^{(i)} \text{trace}[a_i a_i^* G(w)^{-1}] \\ &= \text{trace} \left[ \left( \sum_i w^{(i)} a_i a_i^* \right) G(w)^{-1} \right] \\ &= \text{trace } I = \dim \mathbf{E}^* = n, \end{aligned}$$

where  $I: \mathbf{E}^* \rightarrow \mathbf{E}^*$  is the identity operator. The general case is handled by transforming it to the nonsingular case above. Indeed, let  $\mathcal{X}$  be a subspace of  $\mathbf{E}$  for which  $G(w)$ , viewed as a map from  $\mathcal{X}$  onto  $\text{range } G(w)$ , is invertible and notice that  $\dim \text{range } G(w) = \text{rank } G(w)$ .  $\square$

Let us illustrate the lemma with an example:

**Example 43.** Assume  $G := G(w)$  is of rank 1 and let  $w_1 = 1$ ; all other weights being zero. Since  $G = a_1 a_1^*$ , the solution set of the system  $Gx = a_1$  consists precisely of the vectors  $x$  satisfying  $\langle a_1, x \rangle = 1$ . However,  $\sum_{w^{(i)} > 0} w^{(i)} \gamma_i = \gamma_1 = \langle a_1, x \rangle = 1 = \text{rank } G(w)$ .

The above lemma implies that there is always some index  $i$  such that  $\gamma_i = O(n)$ . However, we already have a procedure for picking  $j$ , and it does not take  $\gamma_j$  into consideration. It would be interesting to see if it is possible to devise a procedure that would guarantee *both* a sufficient decrease in the objective function *and* a reasonable bound on  $\gamma_j$ . Let us remark that the “even better decrease” strategy for choosing  $j$  given in (43) is biased towards choosing one with small  $\gamma_j$ .

Note that, as a corollary of the above lemma, we get the following, albeit somewhat weaker, bound on  $\gamma_j$ :

$$\gamma_j \leq \frac{\text{rank } G(w)}{w^{(j)}}. \tag{54}$$

## 9.4 An alternative proof of the bound on $\gamma_j$

Consider the concave quadratic  $x \mapsto 2\langle a_j, x \rangle - \langle G(w)x, x \rangle$  and observe that its maximizers are precisely the points  $x$  for which  $G(w)x = a_j$ . If  $x_j$  is any such point then

$$\begin{aligned} \gamma_j = \langle a_j, x_j \rangle &= \max_x \{2\langle a_j, x \rangle - \langle G(w)x, x \rangle\} \\ &= \max_x \left\{ 2\langle a_j, x \rangle - \sum_{i=1}^m w^{(i)} \langle a_i, x \rangle^2 \right\} \\ &\leq \max_x \left\{ 2\langle a_j, x \rangle - w^{(j)} \langle a_j, x \rangle^2 \right\} \\ &= \max_{\tau} \left\{ 2\tau - w^{(j)} \tau^2 \right\} = \frac{1}{w^{(j)}}, \end{aligned}$$

yielding another proof of (51). The author wishes to thank Yurii Nesterov for this elegant proof.

## 10 Interpretation

We have seen in Theorem 21 that by solving (resp. approximately solving) problem  $(P3)$ , we have simultaneously solved (resp. approximately solved) also problems  $(P1)$ ,  $(D1)$ ,  $(D'1)$ ,  $(P2)$ ,  $(D2)$  and  $(D3)$ . Moreover, the former theorem mentions how to explicitly construct feasible points for the above problems given a feasible point of  $(P3)$  (with the exception of the last problem). We can therefore, in principle, rewrite our algorithms, which were motivated by problem  $(P3)$ , in terms of iterates feasible for each of the above problems.

For example, if  $\{w_k\}$  is a sequence of iterates produced by Algorithm 2 and  $y_k \in \mathbf{E}$  satisfy  $G(w_k)y_k = d$ , then  $\{v_k\}$  defined by

$$v_k^{(i)} := w_k^{(i)} \langle a_i, y_k \rangle, \quad i = 1, 2, \dots, m,$$

is a sequence of points feasible for  $(D2)$ . Is there a natural way to interpret these iterates in the context of problem  $(D2)$ ?

### 10.1 $(P3)$ : The Frank-Wolfe algorithm on the unit simplex

We will start with an alternative interpretation of our last two algorithms as applied to the main problem of this paper:

$$\boxed{\psi^* := \min_w \{\|d\|_{G(w)}^* : w \in \Delta_m\}.} \quad (P3)$$

The Frank-Wolfe algorithm [11] is a method for solving smooth convex minimization problems over a polytope given as a convex hull of points. At each iteration the objective function is replaced by its linear approximation at the current point. After this, one finds a vertex of the feasible region minimizing the linear approximation — this is a simple enumeration problem. The next iterate is then obtained by performing a line search on the line segment joining the current point and the vertex obtained using the enumeration procedure described above. The

line search can be modified by allowing for Wolfe’s “away steps” [31] — steps in the direction opposite to that towards the vertex maximizing the linear approximation.

It is straightforward to show, using the formula for the derivative of  $\psi^2$  established in Proposition 16, that Algorithm 2 can be interpreted as a Frank-Wolfe method using the former version of line search (the decrease and drop steps of Algorithm 3 correspond to Wolfe’s away steps). Indeed, the linear approximation of  $\psi^2$  at point  $w$  for which  $G(w)$  is invertible is

$$\begin{aligned}\psi^2(w) + D\psi^2(w)(\bar{w} - w) &= \psi^2(w) - \langle G(\bar{w} - w)y, y \rangle \\ &= \psi^2(w) + \langle G(w)y, y \rangle - \langle G(\bar{w})y, y \rangle,\end{aligned}$$

where  $y = G(w)^{-1}d$ . The linearized subproblem can therefore be written as

$$\min_{\bar{w} \in \Delta_m} -\langle G(\bar{w})y, y \rangle = \max_{\bar{w} \in \{e_1, \dots, e_m\}} \sum_i \bar{w}^{(i)} \langle a_i, y \rangle^2.$$

Notice that  $w = e_j$  where  $j = \arg \max_i |\langle a_i, y \rangle|$  solves the above problem. The Frank-Wolfe line search now corresponds to the problem of minimizing  $\psi^2(w(\kappa))$  for  $\kappa \in [0, \infty]$  since  $w(\kappa) = (w + \kappa e_j)/(1 + \kappa)$  parameterizes the line segment joining  $w$  and  $e_j$ . Notice that although in our line search we allow  $-w^{(j)} \leq \kappa < 0$ , the optimal steplength  $\kappa^*$  is always nonnegative (Corollary 30).

For problems where the feasible region is a unit simplex and where the objective function enjoys certain regularity properties such as strong convexity (our function does not satisfy them), it is known that the Frank-Wolfe algorithm with away steps converges linearly [31], [13].

Methods analogous to Algorithm 2 (also interpretable as performing Frank-Wolfe iterations), for computing the minimum volume enclosing ellipsoid of a centrally symmetric body, were proposed by Khachiyan [15], Todd and Yildirim [30]. The method of Todd and Yildirim is a modification of Khachiyan’s algorithm using away steps and has been later analyzed by Ahipařaoglu, Todd and Sun [1] who established its linear convergence. These algorithms, although perhaps without modern convergence analysis, were much earlier independently developed in the statistical community in the context of optimal design by Fedorov [10], Wynn [32], Atwood [3], Silvey [28], and others.

## 10.2 (P2): An ellipsoid-squeezing method for centrally-symmetric LP

Here we will consider problem (P2):

$$\boxed{\frac{1}{\varphi^*} = \max_z \{\langle d, z \rangle : z \in Q^\circ\}.} \tag{P2}$$

Recall that for all  $w \in \Delta_m$  the polar of the ellipsoid  $\mathcal{B}(G(w))$  contains the polar of  $Q$ , that is,  $\mathcal{B}^0(G(w)) \supset Q^0$  (Lemma 6). We also know that

$$\max\{\langle d, y \rangle : y \in \mathcal{B}^0(G(w))\} = \|d\|_{G(w)}^* \geq \psi^* = \frac{1}{\varphi^*}.$$

Let us fix some  $w \in \Delta_m$  and let  $G := G(w)$ . Also let  $y$  be such that  $G(w)y = d$ .

In one iteration of Algorithm 2 (or Algorithm 3) we update  $G$  in a rank-one-and-scaling fashion to  $G(\kappa)$  so as to minimize the value of  $\psi(\kappa) = \|d\|_{G(\kappa)}^*$ . The geometry of this update is rather revealing. Loosely speaking, we choose the step-size parameter  $\kappa$  so as to “push” the polar ellipsoid  $\mathcal{B}^0(G(\kappa))$  by the supporting hyperplane  $\mathcal{H}_d(\kappa) := \{z : \langle d, z \rangle = \|d\|_{G(\kappa)}^*\}$  as far as possible towards  $z^*$ , the optimal point of (P2). This is reminiscent of the correspondence established by Todd and Yildirim [30] between Khachiyan’s ellipsoidal rounding algorithm [15] and the deepest cut ellipsoid method using two-sided symmetric cuts.

Note that  $y/\|d\|_G^*$  lies in the intersection of  $\mathcal{B}(G)$  and  $\mathcal{H}_d := \mathcal{H}_d(0)$  and that  $z := y/\varphi(y)$  is on the boundary of  $Q^0$  and hence is feasible for (P2). This is the current iterate from the perspective of problem (P2). We see our method produces a sequence of points on the boundary of  $Q^0$ .

### 10.3 (D2): An Iteratively Reweighted Euclidean Projection Algorithm

Recall problem (D2):

$$\boxed{\min_v \{\|v\|_1 : Av = d, v \in \mathbf{R}^m\}.} \quad (D2)$$

By Lemma 7, if  $w$  is feasible for (P3) and if  $y \in \mathbf{E}$  is such that  $G(w)y = d$ , then  $v \in \mathbf{R}^m$  defined by  $v^{(i)} := w^{(i)}\langle a_i, y \rangle$ ,  $i = 1, 2, \dots, m$ , is feasible for (D2) and, moreover,

$$\|v\|_1 \leq \|d\|_{G(w)}^*. \quad (55)$$

If we let  $W := \text{Diag}(w)$  (i.e.  $W$  is the diagonal matrix with the entries of vector  $w$  on its diagonal), then, assuming  $G(w)$  is invertible, the above definition of  $v = v(w)$  can be written as

$$v(w) = WA^*y = WA^*(G(w))^{-1}d = WA^*(AWA^*)^{-1}d. \quad (56)$$

We claim that if  $w^{(i)} > 0$  for all  $i$ , which is the case in Algorithm 2, then the point  $v$  above can be obtained as the (unique) minimizer of a certain  $\ell_2$  projection problem. For  $w \in \text{rint } \Delta_m$  and  $W = \text{Diag}(w)$  consider

$$\boxed{\min_v \{\|W^{-1/2}v\|_2 : Av = d, v \in \mathbf{R}^m\}.} \quad (D'2)$$

This problem arises from (D2) if we replace the  $\ell_1$ -norm by the  $\ell_2$ -norm preconditioned by the inverse of a positive definite diagonal matrix with unit trace. Since the set of minimizers does not change if we further replace the objective function by the quadratic  $\frac{1}{2}\|W^{-1/2}v\|_2^2$ , the (necessary and sufficient) KKT conditions for (D'2) are

$$W^{-1}v \in \text{range } A^*, \quad Av = d,$$

from which we readily see that the (unique) minimizer of (D'2) coincides with  $v(w)$  as defined in (56):

$$v^*(w) = WA^*(AWA^*)^{-1}d = v(w).$$

Algorithm 2, as applied to (D2), can therefore be interpreted as follows. At every iteration we maintain a vector of positive weights  $w$  which defines a Euclidean norm on  $\mathbf{R}^m$  by  $v \rightarrow \|W^{-1/2}v\|_2$ . We then “find” the smallest feasible vector in this norm, update the weights and repeat. The weight  $w$  is updated to  $w(\kappa)$  as in (22). As we have discussed before, the arithmetic complexity of every iteration is only  $O(mn)$ , which is the work needed to compute  $A^*x$  for a given vector  $x$ .

### 10.3.1 Two remarks

Let us make two additional observations. First, if we wish to define

$$j := \arg \max_i |\langle a_i, y \rangle|$$

in terms of  $v^*(w)$ , then it is the index for which

$$|[W^{-1}v^*(w)]_j| = \|W^{-1}v^*(w)\|_\infty.$$

Second, notice that

$$\begin{aligned} \|W^{-1/2}v^*(w)\|_2 &= \langle W^{-1/2}WA^*(AWA^*)^{-1}d, W^{-1/2}WA^*(AWA^*)^{-1}d \rangle^{1/2} \\ &= \langle d, (AWA^*)^{-1}AW^{1/2}W^{1/2}A^*(AWA^*)^{-1}d \rangle^{1/2} \\ &= \langle d, (AWA^*)^{-1}d \rangle^{1/2} \\ &= \|d\|_{G(w)}^*, \end{aligned}$$

and hence inequality (55) can be written as

$$\|v^*(w)\|_1 \leq \|W^{-1/2}v^*(w)\|_2.$$

### 10.3.2 The first iterate

Let  $w_0 = (\frac{1}{m}, \dots, \frac{1}{m})$  denote, as usual, the first iterate of Algorithm 2 (resp. Algorithm 3). Then if  $W_0 := \text{Diag}(w_0) = \frac{1}{m}I_m$ , we get

$$v_0 := v(w_0) = W_0A^*(AW_0A^*)^{-1}d = A^*(AA^*)^{-1}d.$$

It is easy to see that this is the shortest feasible vector in the  $\ell_2$  norm. Since  $\|v\|_1 \geq \|v\|_2 \geq \frac{1}{\sqrt{m}}\|v\|_1$  for all  $v \in \mathbf{R}^m$ , then if  $v^*$  is any minimizer of (D2), we have

$$\|v_0\|_1 \leq \sqrt{m}\|v_0\|_2 \leq \sqrt{m}\|v^*\|_2 \leq \sqrt{m}\|v^*\|_1.$$

This shows that the initial iterate  $v_0$  is  $(\sqrt{m} - 1)$ -approximate minimizer of (D2).

## 11 Applications

In this section we apply our algorithm to two problems both of which can be expressed in the form (P3). The first one is the *truss topology design* — a civil engineering application. The second is statistical in nature — the computation of a *c*-optimal design.

### 11.1 Truss topology design

A *truss* is a construction composed of a network of bars linked to one another such as a crane, scaffolding, bridge, wire-model, etc. One can think of a truss as a graph in two or three dimensions. The graph-theoretic terminology then translates as follows: arcs are called bars, vertices are called nodes.

The nodes of a truss are of two categories: *free nodes* and *rigid nodes*. The rigid nodes are attached to some force-absorbing object such as a wall or the ground. Free nodes are subjected to an external force — a *load*. As a result of the load, the free nodes get displaced and bars joining them stretched or squeezed until the structure assumes an equilibrium position in which the internal tensions in the bars compensate for the external forces. A loaded truss therefore stores a certain amount of potential energy called *compliance*. The more there is of this stored energy, the more sensitive the truss becomes to additional loads and/or load variations. It is therefore desirable to design trusses with as small a compliance as possible, given a collection of loads. In this example we will only describe the situation with a fixed vector of loads acting at the free nodes. It is certainly interesting to also consider the case of load scenarios, or perhaps of a dynamic load. These problems are much harder and are out of the scope of this paper.

#### 11.1.1 The problem

The problem we will consider is the following: *Given a set of free and rigid nodes, a set of possible bar locations, a total weight limit on the truss and a vector of external forces acting on the free nodes, design a truss, i.e. give the locations of the bars and their weights, which is capable of holding the given load and has minimum compliance.*

#### 11.1.2 Correspondence with the setting of problem (P3)

The actual derivation of the model can be found, for example, in [4]. Let us describe the parameters of the model in terms of the notation of (P3).

The matrix  $G(w) = \sum_{i=1}^m w^{(i)} a_i a_i^T$  is the *bar-stiffness matrix*. The vector  $w$  corresponds to the weights of the individual tentative bars, normalized so that the total weight of the bars is 1. Let  $p$  be the number of the free nodes. Then we have  $a_i \in \mathbf{R}^n$  where either  $n = 2p$  or  $n = 3p$ , depending on whether we have a  $2d$  or a  $3d$  truss. The system

$$G(w)y = d$$

corresponds to the equilibrium equation between the vector of forces  $d$  acting at the free nodes and the vector of displacements  $y$  of the free nodes. The compliance is one half of the objective

function of problem (P3) squared:

$$Compliance = \frac{1}{2}(\|d\|_{G(w)}^*)^2.$$

We see that problem (P3) corresponds exactly to the truss topology design problem.

### 11.1.3 Three examples

**Example 44.** A unit vertical download force is applied to the right-bottom node of each of the following three  $2d$  trusses:

- (a) A  $3 \times 3$  truss with 3 fixed nodes attached to a wall (the nodes on the left) and 6 free nodes. Hence  $n = 2 \times 6 = 12$ . We allow for tentative bars to be placed among any pair of nodes, with the exception of pairs where there is “overlap” with a chain of other smaller bars. For example, we do not allow placing a bar on the diagonal since this consists of 2 smaller tentative bars already. The number of such potential bars is  $m = 28$ .
- (b) A  $5 \times 5$  truss with 5 fixed nodes. We have  $n = 2 \times 5 \times 4 = 40$  — the number of free nodes times 2. The number of tentative bars is  $m = 400$ .
- (c) A  $9 \times 9$  truss with 9 fixed nodes. In this case  $n = 144$  and  $m = 2040$ .

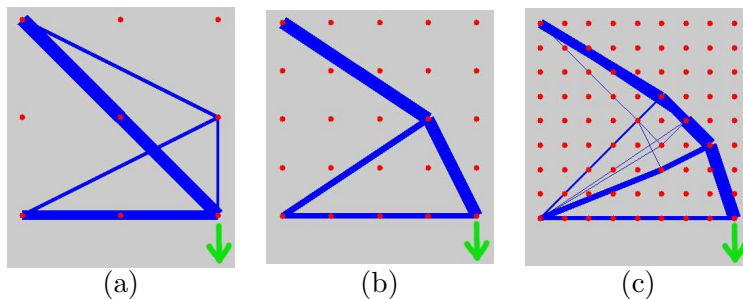


Figure 3: Three optimal trusses.

Figure 3 displays the (approximately) optimal trusses computed with Algorithm **IncDec**. Bars of small weight were removed from the figure. The author wishes to thank Michal Kočvara for sharing his MATLAB code for producing the pictures of the trusses.

Figure 4 lists the performance of our **IncDec** method applied to the three problems of Example 44, with two different accuracy requirements. All computations were done in MATLAB. Let us note that the small  $3 \times 3$  problem was solved by the implementation of the simplex method in MATLAB in 0.5 seconds, the medium  $5 \times 5$  problem in 0.78 seconds, while the large  $9 \times 9$  problem could not be solved by the simplex method within 30 minutes. An interior-point algorithm, however, solved the problem to high accuracy in 0.96 seconds and only 14 iterations.

Truss	$n$	$m$	$\epsilon = 10^{-1}$	$\epsilon = 10^{-4}$	Iteration #	Iteration #
			Time	Time		
$3 \times 3$	12	28	0.07	0.07	413	435
$5 \times 5$	40	200	0.15	1.39	676	7850
$9 \times 9$	144	2040	10.77	367	4450	158601

Figure 4: Performance of Algorithm 3 on three TTD problems.

## 11.2 Optimal design of statistical experiments

The presentation of this subsection is largely based on that in Pukelsheim [23]. See also Fedorov [10] and Silvey [28].

Consider the following situation. An experimenter observes a certain scalar quantity  $y$  which is assumed to depend *linearly* on a vector  $x \in \mathbf{R}^n$  of conditions under his control (a *regression* vector) and a vector of parameters  $\theta \in \mathbf{R}^n$  of interest to him. The observation and/or the model is subject to an additive error  $e$ :

$$y = x^T \theta + e. \quad (57)$$

We will assume that the regression vector  $x$  can be chosen from among a finite collection of vectors  $a_1, \dots, a_m$ , which correspond to the vectors defining  $Q$  — one of the central objects of this paper, known as the *Elfving set* in the statistics community. In what follows we identify  $\mathbf{E}$  and  $\mathbf{E}^*$  with  $\mathbf{R}^n$ .

The statistician wants to estimate a certain function of the parameter  $\theta$  and, in order to do so, decides to observe the outcome under conditions  $x_1, \dots, x_l$ . This is called an *experimental design* of sample size  $l$ . The goal is to construct a design leading to an *unbiased linear estimator*, optimal in a certain sense. Since we restrict the choice of the regression vectors to the finite set  $\{a_1, \dots, a_m\}$ , any design can be described by assigning frequencies to the vectors  $a_i$ . Due to the constraint on the number of observations and the resulting combinatorial structure of feasible frequencies, this approach is usually hard to tackle theoretically. One can instead assign a *weight*  $w^{(i)}$  to each vector  $a_i$ , representing the portion of the entire experiment to be spent under the conditions corresponding to this regression vector.

Let us assume that the errors  $e_j$  are independent random variables with mean zero and (unknown) constant variance  $\sigma^2$  (a *nuisance* parameter). The Fisher information matrix of a design assigning weight  $w^{(i)}$  to point  $a_i$  is given by  $G(w) = \sum_i w^{(i)} a_i a_i^T$ . Now consider the following cases:

- If we wish to minimize the sum of variances of estimators of the individual parameters  $\theta_i$ , this amounts to the problem of minimizing the trace of  $G(w)^{-1}$ . This criterion is referred to as *A-optimality*.
- If the goal is to minimize the variance of the (best unbiased linear estimator) of a linear function of the parameter, say  $c^T \theta$ , it turns out that we need to find  $w \in \Delta_m$  minimizing

$c^T G(w)^{-1} c = (\|c\|_{G(w)}^*)^2$ . If we let  $d := c$ , this is equivalent to our main problem (P3) and is referred to as the  $c$ -optimality criterion in the statistical literature.

- If we wish to minimize the volume of the confidence ellipsoid for  $\theta$ , this corresponds to the problem of maximizing the determinant of  $G(w)$ . This is called the  $D$ -optimality criterion.

Problem (P3) is therefore equivalent to finding the minimum variance unbiased linear estimator of a linear function of the parameter in a statistical linear model with moment assumptions and independent errors.

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

## References

- [1] D. Ahipařaođlu, P. Sun, and M. J. Todd. Linear convergence of a modified Frank-Wolfe algorithm for computing minimum-volume enclosing ellipsoids. Technical Report TR1452, Cornell University, School of Operations Research and Information Engineering, 2006.
- [2] S. Arora, E. Hazan, and S. Kale. The multiplicative weights update method: a meta-algorithm and applications. Technical report, Princeton University, 2005.
- [3] C. L. Atwood. Optimal and efficient design of experiments. *The Annals of Mathematical Statistics*, 40:1570–1602, 1969.
- [4] A. Ben-Tal and A. Nemirovski. *Lectures on Modern Convex Optimization: Analysis, Algorithms, and Engineering Applications*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2001.
- [5] E. Candès, J. Romberg, and T. Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on Information Theory*, 52(2):489–509, 2006.
- [6] E. Candès and T. Tao. Near-optimal signal recovery from random projections and universal encoding strategies. *IEEE Transactions on Information Theory*, 52(12):5406–5425, 2006.
- [7] S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM Journal on Scientific Computing*, 20:33–61, 1999.
- [8] I. Daubechies, R. DeVore, M. Fornasier, and C. S. Güntürk. Iteratively re-weighted least squares minimization for sparse recovery. *Manuscript*, 2008.
- [9] D. L. Donoho. For most large underdetermined systems of linear equations the minimal  $\ell_1$ -norm solution is also the sparsest solution. Technical report, Communications on Pure and Applied Mathematics, 2004.

- [10] V. V. Fedorov. *Theory of Optimal Experiments*. Academic Press, New York, 1972.
- [11] M. Frank and P. Wolfe. An algorithm for quadratic programming. *Naval Research Logistics Quarterly*, 3:95–110, 1956.
- [12] G. H. Golub and Ch. F. Van Loan. *Matrix Computations*. The Johns Hopkins University Press, 1996.
- [13] J. Guélat and P. Marcotte. Some comments on Wolfe’s ‘away step’. *Mathematical Programming*, 35:110–119, 1986.
- [14] J. Hartung. An extension of Sion’s minimax theorem with an application to a method for constrained games. *Pacific Journal of Mathematics*, 103:401–408, 1982.
- [15] L. G. Khachiyan. Rounding of polytopes in the real number model of computation. *Mathematics of Operations Research*, 21:307–320, 1996.
- [16] B. K. Natarajan. Sparse approximate solutions to linear systems. *SIAM J. Comput.*, 24(2):227–234, 1995.
- [17] A. Nemirovski and D. Yudin. *Informational Complexity and Efficient Methods for Solution of Convex Extremal Problems*. J. Wiley and Sons, New York, 1983.
- [18] Yu. Nesterov. A method for unconstrained convex minimization problem with the rate of convergence  $O(\frac{1}{k^2})$ . *Doklady AN SSSR (translated as Soviet. Math. Docl.)*, 269(3):543–547, 1983.
- [19] Yu. Nesterov. *Introductory Lectures on Convex Optimization. A Basic Course*, volume 87 of *Applied Optimization*. Kluwer Academic Publishers, Boston, 2004.
- [20] Yu. Nesterov. Smooth minimization of non-smooth functions. *Mathematical Programming*, 103(1):127–152, 2005.
- [21] Yu. Nesterov. Rounding of convex sets and efficient gradient methods for linear programming problems. *CORE Discussion Paper #2004/04*, January 2004.
- [22] Yu. Nesterov. Unconstrained convex minimization in relative scale. *CORE Discussion Paper #2003/96*, November 2003.
- [23] F. Pukelsheim. *Optimal Design of Experiments (Classics in Applied Mathematics)*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2006.
- [24] P. Richtárik. *Some Algorithms for Large-Scale Linear and Convex Minimization in Relative Scale*. PhD thesis, Cornell University, School of Operations Research and Information Engineering, August 2007.
- [25] P. Richtárik. Approximate level method. *CORE Discussion Paper #2008/83*, December 2008.

- [26] R. T. Rockafellar. *Convex Analysis*. Princeton University Press, Princeton, NJ, USA, 1997. Reprint of the 1970 original, Princeton Paperbacks.
- [27] N. Z. Shor. *Minimization Methods for Nondifferentiable Functions*. Springer-Verlag, Berlin, 1985.
- [28] S. D. Silvey. *Optimal Design: An Introduction to the Theory for Parameter Estimation*. Chapman and Hall, New York, 1980.
- [29] M. Sion. On general minimax theorems. *Pacific Journal of Mathematics*, 8:171–176, 1958.
- [30] M. J. Todd and E. A. Yildırım. On Khachiyan’s algorithm for the computation of minimum volume enclosing ellipsoids. Technical Report TR1435, Cornell University, School of Operations Research and Information Engineering, 2005.
- [31] P. Wolfe. Convergence theory in nonlinear programming. In J. Abadie, editor, *Integer and Nonlinear Programming*, pages 1–36, North-Holland, Amsterdam, 1970.
- [32] H. P. Wynn. The sequential generation of D-optimum experimental design. *The Annals of Mathematical Statistics*, 41:1655–1664, 1970.