

First-order methods for the convex hull membership problem

Rafaela Filippozzi^a, Douglas S. Gonçalves^{a,*}, Luiz-Rafael Santos^b

^aDepartment of Mathematics, Federal University of Santa Catarina, Florianopolis, SC, 88040-900, Brazil

^bDepartment of Mathematics, Federal University of Santa Catarina, Blumenau, SC, 89065-300, Brazil

Abstract

The convex hull membership problem (CHMP) consists in deciding whether a certain point belongs to the convex hull of a finite set of points, a decision problem with important applications in computational geometry and in foundations of linear programming. In this study, we review, compare and analyze first-order methods for CHMP, namely, Frank-Wolfe type methods, Projected Gradient methods and a recently introduced geometric algorithm, called Triangle Algorithm (TA). We discuss the connections between this algorithm and Frank-Wolfe, showing that TA can be interpreted as an inexact Frank-Wolfe. Despite this similarity, TA is strongly based on a theorem of alternatives known as distance duality. By using this theorem, we propose suitable stopping criteria for CHMP to be integrated into Frank-Wolfe type and Projected Gradient, specializing these methods to the membership decision problem. Interestingly, Frank-Wolfe integrated with such stopping criteria coincides with a greedy version of the Triangle Algorithm which is, in its turn, equivalent to an algorithm due to von Neumann. We report numerical experiments on random instances of CHMP, carefully designed to cover different scenarios, that indicate which algorithm is preferable according to the geometry of the convex hull and the relative position of the query point. Concerning potential applications, we present two illustrative examples, one related to linear programming feasibility problems and another related to image classification problems.

Keywords: Convex programming, Convex hull membership problem, Triangle algorithm, Frank-Wolfe algorithms.

1. Introduction

Let $\mathcal{A} := \{v_1, v_2, \dots, v_n\} \subset \mathbb{R}^m$ be given and consider a point $p \in \mathbb{R}^m$. The *convex hull membership problem* (CHMP) consists in deciding whether

$$p \in \text{conv}(\mathcal{A}), \quad (1)$$

where $\text{conv}(\mathcal{A})$ denotes the convex hull of \mathcal{A} . This problem is related to fundamental concepts in linear programming and finds important applications in computational geometry.

Throughout this paper, we denote by $e \in \mathbb{R}^n$ the vector whose n components are all equal to one, by $v^T p$ the Euclidean inner product between vectors $v, p \in \mathbb{R}^m$ and by $\|\cdot\|$ its induced Euclidean norm. We also denote the Euclidean distance between v and p as $d(v, p) := \|v - p\|$, the Euclidean ball centered in v with radius ρ as $B_\rho(v) := \{p \in \mathbb{R}^m \mid d(v, p) < \rho\}$ and the convex combination or segment between v and p as $[v, p]$. The proofs presented are three fold: for new results, when the proof is different from previous proofs in the literature or when the proof contributes to a subsequent discussion.

Let $A := [v_1 \ v_2 \ \dots \ v_n] \in \mathbb{R}^{m \times n}$ be the matrix where each column is one of the n points of \mathcal{A} . One can see that $p \in \text{conv}(\mathcal{A})$ if and only if p is a convex combination of the columns of A . Thus, we can formally describe the convex hull membership problem (1) as the following decision problem:

$$\text{Is there any } x \in \mathbb{R}^n \text{ such that } Ax = p, e^T x = 1, x \geq 0? \quad (2)$$

Problem (2) is a linear programming feasibility problem whose affirmative answer ensures $p \in \text{conv}(\mathcal{A})$. Such feasibility problem can also be cast as the following quadratic programming problem:

$$\min_{x \in \Delta_n} \frac{1}{2} \|Ax - p\|^2 =: \Phi(x), \quad (3)$$

*Corresponding author

Email addresses: rafaela.filippozzi@gmail.com (Rafaela Filippozzi), douglas.goncalves@ufsc.br (Douglas S. Gonçalves), l.r.santos@ufsc.br (Luiz-Rafael Santos)

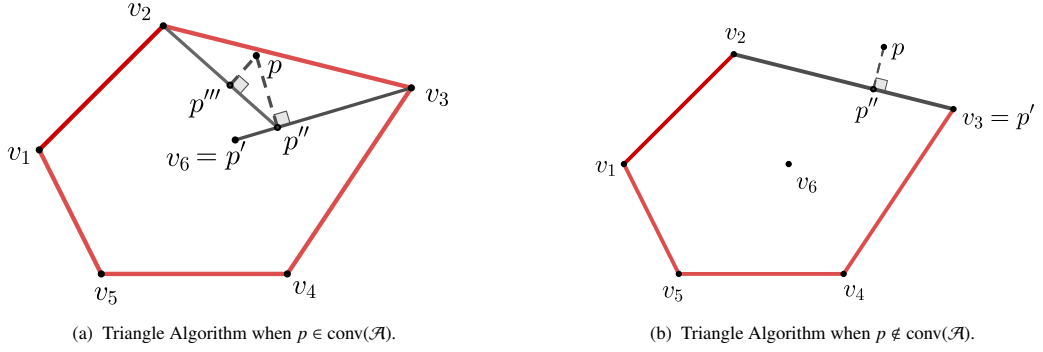


Figure 1: Iterations of the Triangle Algorithm.

where $\Delta_n := \{x \in \mathbb{R}^n \mid e^T x = 1, x \geq 0\}$ is the unit simplex in \mathbb{R}^n . Clearly, $p \in \text{conv}(\mathcal{A})$ when the optimal value of (3) is zero. Furthermore, if x is a feasible point of (3), then $y := Ax \in \text{conv}(\mathcal{A})$. Therefore, another possible formulation of (1) is given by

$$\begin{aligned} \min_{y \in \mathbb{R}^m} \quad & \frac{1}{2} \|y - p\|^2 =: \Psi(y) \\ \text{s.t.} \quad & y \in \text{conv}(\mathcal{A}), \end{aligned} \quad (4)$$

which achieves a zero optimal value if and only if $p \in \text{conv}(\mathcal{A})$.

With formulations (2), (3) and (4) in hand, one could simply apply classic methods for linear or quadratic programming in order to solve CHMP. Nevertheless, Kalantari (2014), inspired by geometric ideas, proposes a specific algorithm for CHMP, called *Triangle Algorithm* (TA), which shows remarkable numerical results (Li & Kalantari, 2013) in comparison with the Simplex method applied to (2) or with Frank-Wolfe (FW) (Frank & Wolfe, 1956) applied to (3).

TA was also successfully employed to solve other important convex hull related problems. For example, Kalantari (2019a) considered the more general problem of separation of two compact convex sets (CHMP is a special case, when one of the sets is a singleton). Awasthi et al. (2020) showed that TA is the building block for an algorithm to enumerate all extreme points of convex hulls in high dimension. For further recent developments involving the Triangle Algorithm see the works of Kalantari (2019b, 2020); Kalantari & Park (2014); Kalantari & Zhang (2022). In particular, Kalantari & Zhang (2022) consider the *Spherical CHMP* where all the points belonging to \mathcal{A} have unit norm and the point p is the origin. The authors showed not only the equivalence between CHMP and spherical CHMP, but also explained how approximate solutions for the former provide approximate solutions for the latter. Furthermore, it was shown that a variant of the Triangle Algorithm for spherical CHMP almost obtains $O(1/\varepsilon)$ complexity as opposed to $O(1/\varepsilon^2)$ complexity of the standard TA and Frank-Wolfe (see Section 3.3 for details).

In general lines, the Triangle Algorithm iterations can be described as follows. For an iterate $p' \in \text{conv}(\mathcal{A})$, the algorithm selects $v_i \in \mathcal{A}$, such that $d(p', v_i) \geq d(p, v_i)$. Such $v_i \in \mathcal{A}$ is called a *pivot*. If a pivot does not exist, we can stop and declare $p \notin \text{conv}(\mathcal{A})$; see Theorem 1.1. Otherwise, update the next iterate to the point lying in the segment $[p', v]$ which is closest to p . The iterations continue until $\|p' - p\| \leq \varepsilon$, in which case it is declared $p \in \text{conv}(\mathcal{A})$. Figure 1 illustrates the scheme (where the sequence of iterates is denoted by p', p'', \dots).

The correctness of the Triangle algorithm follows from the following theorem of alternatives, called *distance duality theorem* and whose proof can be found in Kalantari (2014, Theorem 4).

Theorem 1.1 (Distance duality). *For a given set $\mathcal{A} = \{v_1, v_2, \dots, v_n\} \subset \mathbb{R}^m$ and a point $p \in \mathbb{R}^m$, precisely one of the two conditions is satisfied:*

- (i) *For all $p' \in \text{conv}(\mathcal{A})$, there exists $v_i \in \mathcal{A}$ such that $d(p', v_i) \geq d(p, v_i)$;*
- (ii) *There exists $p' \in \text{conv}(\mathcal{A})$ such that $d(p', v_i) < d(p, v_i)$, for all $i = 1, \dots, n$.*

Note that Theorem 1.1 implies that either for every $p' \in \text{conv}(\mathcal{A})$ there exists $v \in \mathcal{A}$ such that v is closer to p than to p' , or there exists $p' \in \text{conv}(\mathcal{A})$ such that the bisector hyperplane orthogonal to the segment $[p', p]$ separates p from $\text{conv}(\mathcal{A})$ (see Figure 1 in Kalantari (2014)).

Kalantari (2014) proved that the first condition is valid if and only if $p \in \text{conv}(\mathcal{A})$, and the second condition if and only if $p \notin \text{conv}(\mathcal{A})$. If the latter holds, we say that p' is a *witness* that $p \notin \text{conv}(\mathcal{A})$ because the hyperplane

$$H_{[p,p']} := \{y \in \mathbb{R}^m \mid (p - p')^T y = (p - p')^T (p' + p)/2\}$$

strictly separates p from $\text{conv}(\mathcal{A})$. Furthermore, when $p \notin \text{conv}(\mathcal{A})$, for a p -witness p' , we have

$$\frac{1}{2} \|p' - p\| \leq \Delta \leq \|p' - p\|, \quad (5)$$

where $\Delta := \min\{d(v, p) \mid v \in \text{conv}(\mathcal{A})\}$. Thus, the distance from p to $\text{conv}(\mathcal{A})$ is approximated within a factor of two.

The distance duality theorem provides unique features to the Triangle algorithm when compared to first-order algorithms applied to the optimization formulations of CHMP. For instance, one can consider applying the classic FW to (3) or (4). We shall show later (see Section 3.3) that the subproblem that must be solved at each FW iteration corresponds to find $v_i \in \arg \min\{(p' - p)^T v_j \mid v_j \in \mathcal{A}\}$ whereas $v_i \in \mathcal{A}$ is a pivot for TA if and only if $(p' - p)^T v_i \leq (\|p'\|^2 - \|p\|^2)/2$ (see Lemma 2.2). Thus, the first alternative of Theorem 1.1 relaxes this requirement of the “optimal” pivot chosen by FW. Moreover, the second alternative of Theorem 1.1 provides a smart stopping criterion, allowing us to state that $p \notin \text{conv}(\mathcal{A})$ without the need of solving (3) (or (4)) to ε -optimality.

Yet, TA is similar to FW applied to (4) in the sense that it keeps its iterates as convex combinations of points of \mathcal{A} , adding at most one new point per iteration to this combination. What is more, we will show that the so-called pivot of the Triangle algorithm can be seen as an inexact solution to the FW subproblem. Not surprisingly, convergence and iteration complexity results for the Triangle algorithm (and its variants) follow closely those of FW methods: namely, in general the convergence is usually sublinear, and linear convergence is only achieved under stronger assumptions.

The contributions of this paper are the following.

- (i) We study the similarities between Frank-Wolfe and the Triangle algorithm: previous works on the Triangle algorithm (for instance, Kalantari, 2014, 2019a) briefly mention FW for CHMP and contrast its iteration mechanism with TA. Here, we show that FW with a stopping criterion based on distance duality (useful in case $p \notin \text{conv}(\mathcal{A})$) can be interpreted as a greedy version of the Triangle algorithm: such version is equivalent to an algorithm attributed to von Neumann (as communicated by Dantzig, 1992). Moreover, we also show that a pivot from TA can be interpreted as an inexact solution of Frank-Wolfe subproblem, in the sense of (Jaggi, 2013).
- (ii) Based on this relationship, we propose variants of Frank-Wolfe that take advantage of the distance duality.
- (iii) Relying on distance duality, we devise suitable stopping criteria for the considered first-order methods employed to solve CHMP.
- (iv) We compare the performance of the new variants with their classic versions and other first-order methods such as projected gradient methods through comprehensive numerical experiments on CHMP: even though previous works (Awasthi et al., 2020; Li & Kalantari, 2013) have reported some numerical experiments comparing the Triangle Algorithm with Frank-Wolfe, applied to the optimization problem (3), specialized stopping criteria for CHMP were not used for FW. Aiming at a fair comparison, we integrate FW (and the other first-order methods) with specialized stopping criteria. What is more, as far as we know, none of these works have considered Away-Step Frank-Wolfe (or Frank-Wolfe) applied to (4) and Projected Gradient methods applied to (3) in the computational experiments.
- (v) We illustrate the usefulness of these algorithms in two potential applications: linear programming feasibility problems and image classification problems.

For that, this paper is organized as follows. In Section 2 we provide a formal description of the Triangle algorithm along with its main convergence and iteration complexity results. Section 3 is devoted to a brief review of the Frank-Wolfe method for minimization of convex smooth functions over convex compact sets and a variant known as Away-Step Frank-Wolfe. In Section 3.3, we discuss similarities between Triangle and Frank-Wolfe algorithms applied to CHMP. In Section 3.4, by using distance duality in FW we arrive at a greedy variant

of the Triangle Algorithm. Section 4 reviews the spectral projected gradient method and discuss its application to CHMP. In Section 5 we devise specific stopping criteria for Frank-Wolfe and projected gradient methods when applied to the CHMP. The numerical experiments reported in Section 6 show that the new variants of FW and Triangle algorithms perform better than their classical counterparts. We report that in randomly generated problems covering different scenarios (as described in Section 6.1) with respect to the geometry of the convex hull and the relative position of the query point. Two potential applications of these algorithms, one in the linear feasibility problem and the other on image classification problem, are detailed in Section 6.2 and Section 6.3, respectively. Conclusions are drawn in Section 7.

2. Triangle Algorithm

In this section, we formally describe the Triangle algorithm and review the main convergence results from Kalantari (2014). We also provide complementary results that shall be useful later.

The Triangle algorithm takes as input a finite set of points $\mathcal{A} = \{v_1, \dots, v_n\} \subset \mathbb{R}^m$, a point $p \in \mathbb{R}^m$, and a tolerance $\varepsilon \in (0, 1)$. At the k -th iteration, for a given point $p_k \in \text{conv}(\mathcal{A})$, the algorithm looks for $v \in \mathcal{A}$ that satisfies $d(v, p) \leq d(v, p_k)$. This v is called a *pivot*. More specifically, we say that v is a pivot at p_k . By defining $R := \max\{d(v_i, p) \mid v_i \in \mathcal{A}\}$, if $d(p_k, p) \leq \varepsilon R$, we say that p_k is an ε -solution, and the algorithm stops, stating $p \in \text{conv}(\mathcal{A})$. Otherwise, the next iterate p_{k+1} is the closest point to p on the line segment $[p_k, v]$. If there exists no pivot at $p_k \in \text{conv}(\mathcal{A})$, that is, when $d(v_j, p) > d(v_j, p_k), \forall j = 1, \dots, n$, we say that p_k is a *witness* that $p \notin \text{conv}(\mathcal{A})$. In this case, the orthogonal bisecting hyperplane to the line segment $[p, p_k]$ separates p from $\text{conv}(\mathcal{A})$. This iterative scheme is summarized in Algorithm 1.

Algorithm 1: TRIANGLE ALGORITHM (TA)

Data: $\mathcal{A}, p, R, \varepsilon \in (0, 1)$

- 1 Choose $p_0 \in \arg \min\{d(v_j, p) \mid v_j \in \mathcal{A}\}$
- 2 **for** $k = 0, 1, 2, \dots$ **do**
- 3 **if** $d(p_k, p) < \varepsilon R$ **then** stop
- 4 **if** $d(v_j, p) > d(v_j, p_k), \forall j = 1, \dots, n$ **then** stop: p_k is a p -witness.
- 5 Choose $v_j \neq p_k$, such that $d(v_j, p) \leq d(v_j, p_k)$.
- 6 Set $\tilde{\gamma}_k \in \arg \min\{d(p, (1 - \gamma)p_k + \gamma v_j) \mid 0 \leq \gamma \leq 1\}$.
- 7 $p_{k+1} \leftarrow (1 - \tilde{\gamma}_k)p_k + \tilde{\gamma}_k v_j$.
- 8 **end**

Remark 2.1. (i) In step 1 of Algorithm 1 we could have chosen any $p' \in \text{conv}(\mathcal{A})$ as p_0 , however starting with a point of \mathcal{A} closest to p has technical reasons that will become clear ahead.

(ii) If a p -witness is detected at step 4, according to Theorem 1.1, we have that $p \notin \text{conv}(\mathcal{A})$.

(iii) If the Triangle Algorithm stops at step 3, we have a point $p_k \in \text{conv}(\mathcal{A})$ whose distance is less than εR from p , thus we classify p as an element of $\text{conv}(\mathcal{A})$.

The next lemma will be useful in analyzing the variations of the Triangle Algorithm. It features equivalent characterizations for a pivot that follow from simple algebraic manipulation.

Lemma 2.2. *Let $p_k \in \text{conv}(\mathcal{A})$, $v_j \in \mathcal{A}$ and $p \in \mathbb{R}^m$ be given. The following are equivalent:*

- (i) $\|v_j - p\| \leq \|v_j - p_k\|$;
- (ii) $2v_j^T(p_k - p) \leq \|p_k\|^2 - \|p\|^2$;
- (iii) $(p_k - p)^T(v_j - p_k) \leq -\frac{1}{2} \|p_k - p\|^2$;
- (iv) $(p_k - p)^T(v_j - p) \leq \frac{1}{2} \|p_k - p\|^2$.

Using this lemma we can establish the following auxiliary result.

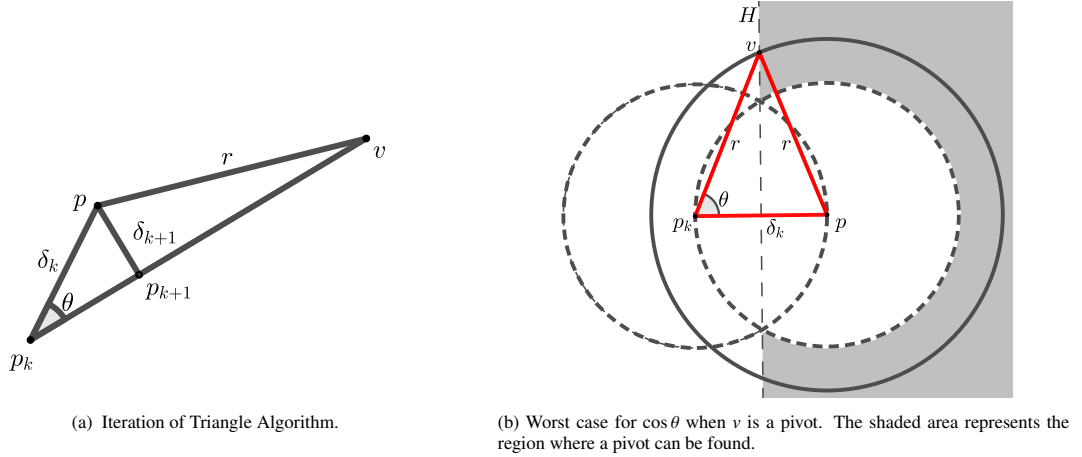


Figure 2: Motivation for the proof of Theorem 2.4.

Lemma 2.3. *In every iteration k of Algorithm 1, if $d(p_k, p) > \varepsilon R$ and if there exists a pivot v_j , then $\bar{\gamma}_k$ from step 6 has a closed formula given by*

$$\bar{\gamma}_k = -\frac{(p_k - p)^T (v_j - p_k)}{\|v_j - p_k\|^2} \in (0, 1]. \quad (6)$$

Moreover, $d(p_{k+1}, p) < d(p_k, p)$.

Proof. It follows from simple calculations that the unconstrained minimizer of the convex quadratic function $\phi(\gamma) := d(p, (1-\gamma)p_k + \gamma v_j)^2$ required in step 6 is given by (6). Let us show that $\bar{\gamma}_k \in (0, 1]$. From Lemma 2.2(iii), we have

$$\bar{\gamma}_k \geq \frac{1}{2} \frac{\|p_k - p\|^2}{\|v_j - p_k\|^2} > 0,$$

where the last inequality follows from the hypothesis $d(p_k, p) > \varepsilon R > 0$. On the other hand, from Cauchy-Schwarz inequality, we obtain

$$\bar{\gamma}_k \leq \frac{\|p_k - p\|}{\|v_j - p_k\|} \leq \frac{\|p_k - p\|}{\|v_j - p\|} \leq 1,$$

where the second inequality follows from step 5 and Lemma 2.2(i). The third inequality is ensured in the first iteration by step 1. Its validity in the subsequent iterations comes from $d(p_{k+1}, p) < d(p_k, p)$ that we prove now. Note that

$$\begin{aligned} d(p_{k+1}, p)^2 &= \|p_{k+1} - p\|^2 = \|(1 - \bar{\gamma}_k)p_k + \bar{\gamma}_k v_j - p\|^2 = \|\bar{\gamma}_k(v_j - p_k) + (p_k - p)\|^2 \\ &= \bar{\gamma}_k^2 \|v_j - p_k\|^2 + \bar{\gamma}_k 2(p_k - p)^T (v_j - p_k) + \|p_k - p\|^2 < \|p_k - p\|^2 = d(p_k, p)^2, \end{aligned} \quad (7)$$

where the last inequality follows from the fact that $\bar{\gamma}_k > 0$ from (6) is the strict unconstrained minimizer of the quadratic in (7). Therefore, $d(p_{k+1}, p) < d(p_k, p)$ and $0 < \bar{\gamma}_k \leq 1$ in every iteration such that $d(p_k, p) > \varepsilon R$ and a pivot v_j exists. \square

As a consequence of Lemma 2.3, we have that the sequence of distances given by $\delta_k := d(p_k, p)$ is monotonically decreasing. We can indeed quantify such reduction by observing that the inequalities $d(p_k, p) \leq d(v, p) \leq d(v, p_k)$ imply that the points p_k, p and v are non-collinear and thus result in a non-degenerate triangle, as depicted in Figure 2a, where the angle $\theta_k := \angle p p_k v \in [0, \pi/2)$. Hence, we have

$$\|p_{k+1} - p\|^2 = (1 - \cos^2 \theta_k) \|p_k - p\|^2. \quad (8)$$

Now, in view of the pivot characterization (Lemma 2.2), we can see that if $\|v - p\| = r$, $\cos \theta_k$ is minimum when v is over the hyperplane $H := \{y \in \mathbb{R}^m \mid (p - p_k)^T y = (\|p\|^2 - \|p_k\|^2)/2\}$. In this case, $\cos \theta_k = \delta_k/2r$; see an illustration in Figure 2b.

From these remarks and from (8) we have the next theorem.

Theorem 2.4 (Kalantari (2014, Theorem 8)). *Let p, p', v be distinct points in \mathbb{R}^m and suppose that $d(v, p) \leq d(v, p')$. Let p'' be the point in the segment $[p', v]$ that is closest to p . Define $\delta := d(p', p)$, $\delta' := d(p'', p)$, $r := d(v, p)$ and assume $\delta \leq r$. Then,*

$$\delta' \leq \delta \sqrt{1 - \frac{\delta^2}{4r^2}}.$$

Remark 2.5. Regarding Algorithm 1, we have

$$\delta_{k+1} \leq \delta_k \sqrt{1 - \frac{\delta_k^2}{4R^2}} \leq \delta_k \exp\left(-\frac{\delta_k^2}{8R^2}\right),$$

where the last inequality follows from $1 + x \leq \exp(x)$; recall that $R = \max\{d(v_j, p) \mid v_j \in \mathcal{A}\}$. We also observe that, as long as $\delta_k > \varepsilon R$, it holds that $\exp\left(-\frac{\delta_k^2}{8R^2}\right) < \exp\left(-\frac{\varepsilon^2}{8}\right) < 1$.

Now, we can formally aggregate the complexity bound of Algorithm 1 in the next result.

Theorem 2.6 (Complexity of TA (Kalantari, 2014, Theorem 9)). *Algorithm 1 correctly solves the convex hull membership problem (1) with the following complexity:*

- (i) *If $p \in \text{conv}(\mathcal{A})$, given $\varepsilon > 0$, $p_0 \in \text{conv}(\mathcal{A})$, with $\delta_0 = d(p, p_0) \leq \min\{d(v_i, p) \mid i = 1, \dots, n\}$, the maximum number of iterations K_ε to compute a point $p_\varepsilon \in \text{conv}(\mathcal{A})$ such that $d(p_\varepsilon, p) < \varepsilon R$ satisfies*

$$K_\varepsilon \leq \frac{48}{\varepsilon^2} = O(\varepsilon^{-2}).$$

- (ii) *If $p \notin \text{conv}(\mathcal{A})$ the number of iterations K_Δ to compute a p -witness is such that*

$$K_\Delta \leq \frac{48R^2}{\Delta^2} = O\left(\frac{R^2}{\Delta^2}\right), \quad (9)$$

where $\Delta = \min\{d(x, p) \mid x \in \text{conv}(\mathcal{A})\}$.

We point out that when $p \notin \text{conv}(\mathcal{A})$ Theorem 2.6(ii) states the iteration complexity of TA depending only on the geometry of the problem (namely, the constants Δ and R). On the other hand, the $O(\varepsilon^{-2})$ iteration complexity in Theorem 2.6(i) implies that the convergence rate is only sublinear. In order to improve this complexity result, a more restrictive definition of pivot is presented next.

Definition 2.7 (Strict pivot (Kalantari, 2014, Definition 8)). *Let $p \in \mathbb{R}^m$ and $p_k \in \text{conv}(\mathcal{A})$. We say that $v \in \mathcal{A}$ is a strict pivot, if the angle $\angle p_k p v$ between the segments $[p_k, p]$ and $[p, v]$ is such that $\angle p_k p v \geq \pi/2$.*

Remark 2.8. Clearly, $v \in \mathcal{A}$ is strict pivot for p_k whenever $(p_k - p)^T(v - p) \leq 0$, which is a stronger condition than the one in Lemma 2.2(iv). We depict the geometric interpretation of a strict pivot in Figure 3, where the hyperplane \bar{H} is defined as $\bar{H} := \{y \in \mathbb{R}^m \mid (p - p_k)^T(y - p) = 0\}$. Analogously to Remark 2.5, in the case of v being a strict pivot, we have

$$\delta_{k+1} = \frac{\delta_k r}{\sqrt{r^2 + \delta_k^2}} \leq \frac{\delta_k R}{\sqrt{R^2 + \delta_k^2}} \leq \delta_k \sqrt{1 - \frac{\delta_k^2}{2R^2}} \leq \delta_k \exp\left(-\frac{\delta_k^2}{4R^2}\right).$$

Note that a distance duality theorem for strict pivots is also available and is presented in sequel.

Theorem 2.9 (Distance duality for a strict pivot (Kalantari, 2014, Theorem 10)). *Assume $p \notin \mathcal{A}$. Then, we have $p \in \text{conv}(\mathcal{A})$ if, and only if, for each $p_k \in \text{conv}(\mathcal{A})$ there exists a strict pivot $v \in \mathcal{A}$.*

Similarly to Theorem 2.6, the Triangle algorithm using strict pivots also requires $O(1/\varepsilon^2)$ iterations to find an ε -solution (albeit the multiplicative constant is smaller; see Kalantari (2014, Theorem 11)), whenever p belongs to $\text{conv}(\mathcal{A})$. Nonetheless, strict pivots shall be useful ahead to showcase a situation where the convergence rate is linear. First, note from (8) that the reduction in $\delta_k = d(p_k, p)$, at each iteration of the Triangle Algorithm, depends on the angle $\theta_k = \angle p p_k v$ (see Figure 2). Thus, if $\cos \theta_k$ stays bounded away from zero, it is possible to improve the complexity results for TA. Indeed, we formalize this fact, which was established in Kalantari (2014), with the aid of the next assumption.

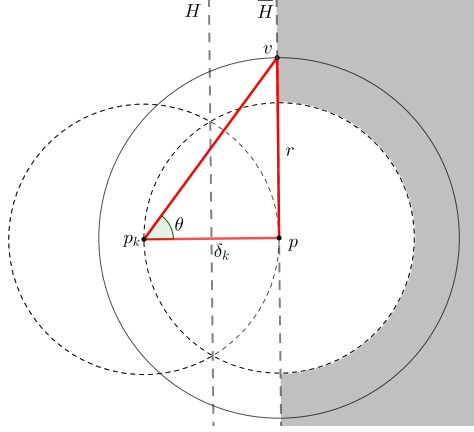


Figure 3: The shaded area represents the region where strict pivots can be found.

Assumption 2.10. There exists a constant $c > 0$ (called *visibility factor*) such that the following inequality is fulfilled

$$\sin \theta_k \leq \frac{1}{\sqrt{1+c}}, \quad \forall p_k \in \text{conv}(\mathcal{A}), \quad (10)$$

where $\theta_k = \angle pp_kv$ and v is a pivot at p_k .

Theorem 2.11 (Kalantari (2014, Theorem 13)). Let $\delta_0 = d(p_0, p)$, $p_0 \in \text{conv}(\mathcal{A})$ and assume that Assumption 2.10 holds.

- (i) If $p \in \text{conv}(\mathcal{A})$, then the number of iterations of the Triangle Algorithm to obtain $p_\varepsilon \in \text{conv}(\mathcal{A})$ such that $d(p_\varepsilon, p) \leq \varepsilon R$ is

$$O\left(\frac{1}{c} \ln \frac{\delta_0}{\varepsilon R}\right).$$

- (ii) If $p \notin \text{conv}(\mathcal{A})$ the number of iterations of the Triangle Algorithm to obtain a p -witness is

$$O\left(\frac{1}{c} \ln \frac{\delta_0}{\Delta}\right).$$

Iteration complexities of Theorem 2.11 are improvements in comparison with those of Theorem 2.6, as long as the visibility factor c is bounded away from zero. As we shall see ahead, this is the case, for example, when p belongs to the relative interior of $\text{conv}(\mathcal{A})$. Nevertheless, in other cases, c can be very close to 0 and, as a consequence, the complexities of Theorem 2.11 may be way worse than those of Theorem 2.6. When this happens, TA may experience a zigzagging phenomenon as illustrated by the following example.

Example 2.12. Consider $\mathcal{A} = \{v_0, v_1, v_2, v_3, v_4\} \subset \mathbb{R}^2$ where $\{v_0, v_1, v_2, v_3\}$ are vertices of the unit square and v_4 is an interior point. Assume that the point of \mathcal{A} closest to the query point p is v_4 (which was generated by shifting the center of the square a bit to the right). We analyze two instances depicted in Figure 4. In Figure 4a, we have the case where the query point p is the midpoint of an edge of $\text{conv}(\mathcal{A})$ whereas Figure 4b shows the case where p is moved to a position outside the square. In both cases, the zigzagging behavior of TA shows up. We only plot the first 60 iterations of TA, in each picture: it takes more than a million iterations to achieve $\|p_k - p\| \leq \varepsilon R$, with $\varepsilon = 10^{-4}$ in case (a) and 81 iterations to find a witness in case (b). For $p \in \text{conv}(\mathcal{A})$, we have $\delta_0 = 0.40$ and $R = 1.12$. We monitored the sequence $(\sin \theta_k)_{k \in \mathbb{N}}$ based on TA iterations in order to estimate an upper bound for (10). We notice that the values approach 1 as k grows; thus, for the points we computed, if the visibility factor c exists, it is not greater than 10^{-7} . When $p \notin \text{conv}(\mathcal{A})$, we get $\delta_0 = 0.45$ and $\Delta = 0.05$. Again, tracking $(\sin \theta_k)_{k \in \mathbb{N}}$ allowed us to provide an upper bound for the visibility constant c of 0.0044.

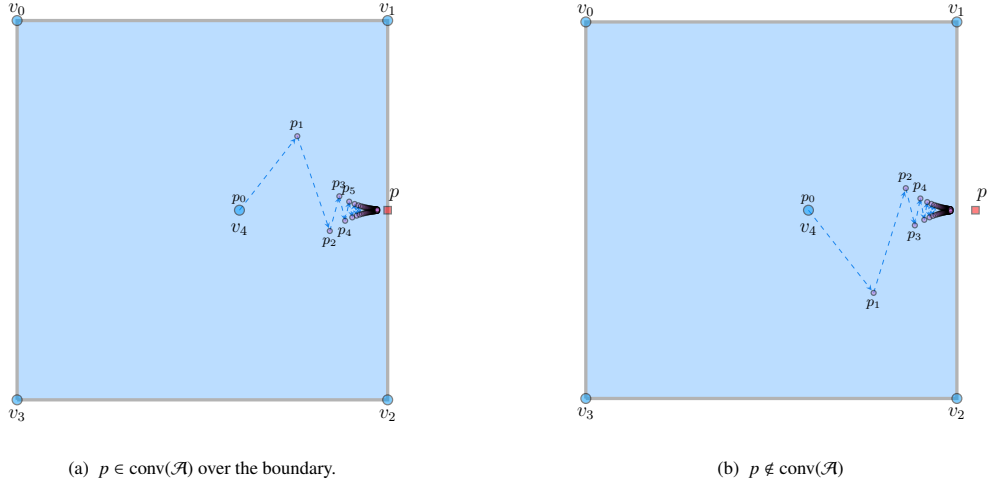


Figure 4: Unit squares of Example 2.12 with the first 60 TA iterates.

Under the additional assumption that p lies in the relative interior of $\text{conv}(\mathcal{A})$, and under a suitable choice of pivots in each iteration, it is possible to show that (10) is satisfied with a constant c whose lower bound depends only on the geometry of the problem.

Assumption 2.13. $B_\rho(p)$ is contained in the relative interior of $\text{conv}(\mathcal{A})$.

Let p_k be the current iterate such that $d(p, p_k) > \varepsilon R$. If Assumption 2.13 holds for a given $\rho > 0$, then there exists a $q \in \text{conv}(\mathcal{A})$ that belongs to the extension of the segment $[p_k, p]$ (see Figure 5a) with $d(q, p) = \rho$. Now, if there exists a strict pivot v with the property $(p - q)^T(v - q) \leq 0$, then

$$\begin{aligned} \rho R \cos \angle qpv &\geq \|p - q\| \|v - p\| \cos \angle qpv = (q - p)^T(v - p) \\ &= (q - p)^T(v - q + q - p) = \|q - p\|^2 - (p - q)^T(v - q) \geq \rho^2 \end{aligned}$$

which implies that

$$\sin \theta_k \leq \sin \angle qpv \leq \sqrt{1 - \rho^2/R^2} \leq \frac{1}{\sqrt{1 + \rho^2/R^2}}. \quad (11)$$

We point out that (11) is necessary in the proof of the next theorem (Kalantari, 2014, Theorem 14), and thus we consider the following assumption.

Assumption 2.14. Let $q = p + \tau(p - p_k)$ for $\tau > 0$ such that $q \in \text{conv}(\mathcal{A})$ and $d(q, p) = \rho$. The Triangle Algorithm uses, at each iteration, a strict pivot v satisfying

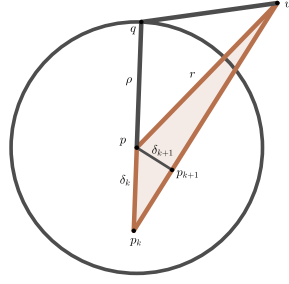
$$(p - q)^T(v - q) \leq 0. \quad (12)$$

Remark 2.15. A strict pivot v satisfying (12) is represented in Figure 5a. We remark that a strict pivot does not necessarily satisfy condition (12). In Figure 5b, the gray area represents the region where a strict pivot can be found. Note that v_1, v_2 and v_3 are strict pivots, but v_2 does not satisfy (12) because it is in the *wrong* side of the hyperplane $\tilde{H}_q := \{y \in \mathbb{R}^m \mid (p - q)^T(y - q) = 0\}$. In fact, only strict pivots belonging to the dashed area satisfy the inequality (12).

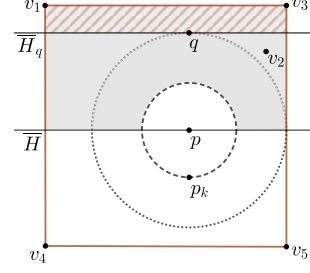
Theorem 2.16. Let $R = \max\{d(v_i, v_j) \mid v_i, v_j \in \mathcal{A}\}$, $p_0 \in \text{conv}(\mathcal{A})$, $\delta_0 = d(p_0, p)$ and suppose that Assumptions 2.13 and 2.14 hold. For any given $\varepsilon \in (0, 1)$, the number of iterations that Algorithm 1 takes to compute $p_\varepsilon \in \text{conv}(\mathcal{A})$, such that $d(p, p_\varepsilon) \leq \varepsilon R$, is

$$O\left(\frac{R^2}{\rho^2} \ln \frac{\delta_0}{\varepsilon R}\right).$$

Proof. It can be found in (Kalantari, 2014, Theorem 14). Even though Assumption 2.14 is not explicit in the statement of Kalantari's theorem, it becomes necessary in the proof. \square



(a) A strict pivot v that satisfies condition (12).



(b) Shaded area is where strict pivots can be found while in the dashed area are the strict pivots for which (12) holds.

Figure 5: Representations of strict pivots.

This result implies the linear convergence of Algorithm 1, using strict pivots satisfying Assumption 2.14, when $p \in B_\rho(p) \subset \text{relint}(\text{conv}(\mathcal{A}))$. In Section 3.4 we shall see a variant of Algorithm 1 for which Assumption 2.14 holds in every iteration.

3. Frank-Wolfe for CHMP

We kick-start this section presenting the Frank-Wolfe algorithm (FW) for convex constrained optimization. We then explore one of its variants called Away Step Frank-Wolfe (ASFW) for which a linear convergence rate can be obtained. Finally, we formally discuss connections between FW and the Triangle Algorithm when applied to the convex hull membership problem.

3.1. Frank-Wolfe algorithm

Let $f : \mathbb{R}^m \rightarrow \mathbb{R}$ be a continuously differentiable convex function whose gradient is L -Lipschitz and $C \subset \mathbb{R}^m$ a non-empty, convex and compact set. Following Lacoste-Julien & Jaggi (2015), we assume that $C := \text{conv}(\mathcal{A})$, where $\mathcal{A} \subset \mathbb{R}^m$ is a finite set of points. Consider the following optimization problem

$$\min_{x \in C} f(x). \quad (13)$$

The Frank-Wolfe algorithm (Frank & Wolfe, 1956) (also known as Conditional Gradient) can be used to solve problem (13), and it is summarized in Algorithm 2. FW is a first-order iterative method for smooth convex optimization which aims to solve (13) through a sequence of linearized subproblems.

Algorithm 2: Frank-Wolfe (FW)

Data: $x_0 \in \mathcal{A}, \epsilon > 0$

- 1 **for** $k = 0, 1, 2, \dots$ **do**
- 2 Set $\bar{x}_k \in \arg \min_{x \in \mathcal{A}} \nabla f(x_k)^T (x - x_k)$
- 3 $d_k \leftarrow \bar{x}_k - x_k$.
- 4 **if** $\nabla f(x_k)^T d_k \geq -\epsilon$ **then stop**
- 5 Choose $\gamma_k \in (0, 1]$.
- 6 $x_{k+1} \leftarrow x_k + \gamma_k d_k$
- 7 **end**

Remark 3.1. One of the main advantages of FW is that, in many applications, the subproblem of Step 2 either admits a closed form solution or is much cheaper in terms of computational cost than the quadratic subproblems required by projected gradient or proximal methods, leading to good scalability (Jaggi, 2013). Moreover, it keeps the iterates as convex combination of few points of \mathcal{A} , which is an interesting feature in problems requiring sparse solutions (Clarkson, 2010).

Remark 3.2. Let x^* be a solution of (13). Since f is convex and differentiable, and by the definition of Step 2, we have

$$h(x_k) := f(x_k) - f(x^*) \leq -\nabla f(x_k)^T (x^* - x_k) \leq -\min_{x \in \mathcal{A}} \nabla f(x_k)^T (x - x_k) =: g(x_k),$$

that is, $g(x_k)$ provides an upper bound for the primal gap $h(x_k)$, therefore justifying the stopping criterion in Step 4.

In the following, we briefly review the convergence analysis for Algorithm 2, so that the comparison of its convergence rate with those of other algorithms can be clear. For further details, see (Beck & Shtern, 2017; Jaggi, 2013; Lacoste-Julien & Jaggi, 2015).

Since ∇f is Lipschitz with constant $L > 0$, we get that

$$f(x_{k+1}) = f(x_k + \gamma_k d_k) \leq f(x_k) + \gamma_k \nabla f(x_k)^T d_k + \frac{L}{2} \gamma_k^2 \|d_k\|^2 \leq f(x_k) - \gamma_k g(x_k) + \gamma_k^2 \frac{L}{2} D^2,$$

where $D := \max\{\|x - y\| \mid x, y \in C\}$ is the diameter of C . From this inequality and the definition of $h(x_k)$, we obtain

$$h(x_{k+1}) \leq (1 - \gamma_k)h(x_k) + \gamma_k^2 \frac{L}{2} D^2.$$

Then, it is possible to derive (see Theorem 1 in Jaggi, 2013) that, for $\gamma_k = \frac{2}{k+2}$,

$$h(x_k) \leq \frac{2LD^2}{k+2}, \quad (14)$$

from which the sublinear convergence of FW is established¹. Unfortunately, this $O(1/k)$ complexity is tight, as showed by Canon & Cullum (1968).

Nevertheless, under strong convexity of f , it is possible to achieve linear convergence in certain cases. In fact, if we suppose the existence of $\mu > 0$ such that, for all x, y ,

$$f(y) \geq f(x) + \nabla f(x)^T (y - x) + \frac{\mu}{2} \|y - x\|^2,$$

and define $e_k := x^* - x_k$, we have

$$f(x_k + \lambda e_k) - f(x_k) \geq \lambda \nabla f(x_k)^T e_k + \frac{\lambda^2}{2} \mu \|e_k\|^2 \geq -\frac{(\nabla f(x_k)^T \hat{e}_k)^2}{2\mu}, \quad \forall \lambda \in \mathbb{R},$$

where $\hat{e}_k := e_k / \|e_k\|$. In particular, for $\lambda = 1$, we get

$$-h(x_k) \geq -\frac{(\nabla f(x_k)^T \hat{e}_k)^2}{2\mu}. \quad (15)$$

On the other hand, using again that ∇f is L -Lipschitz we obtain

$$f(x_k + \gamma_k d_k) - f(x_k) \leq \gamma_k \nabla f(x_k)^T d_k + \frac{\gamma_k^2}{2} L \|d_k\|^2.$$

By taking $\gamma_k := \arg \min_{\gamma \in (0, \gamma_{\max}] } f(x_k + \gamma d_k)$ in Algorithm 2, where $\gamma_{\max} \in (0, 1]$ is the maximum allowed step-size,

if $\gamma_k^* := -\frac{\nabla f(x_k)^T d_k}{L \|d_k\|^2} \leq \gamma_{\max}$, then $\gamma_k = \gamma_k^*$ and

$$\gamma_k \nabla f(x_k)^T d_k + \frac{\gamma_k^2}{2} L \|d_k\|^2 = -\frac{(\nabla f(x_k)^T \hat{d}_k)^2}{2L},$$

¹If $\tilde{\gamma}_k$ is the exact line-search step-size, the previous analysis is still valid because $f(x_k + \tilde{\gamma}_k d_k) \leq f(x_k + \frac{2}{k+2} d_k)$.

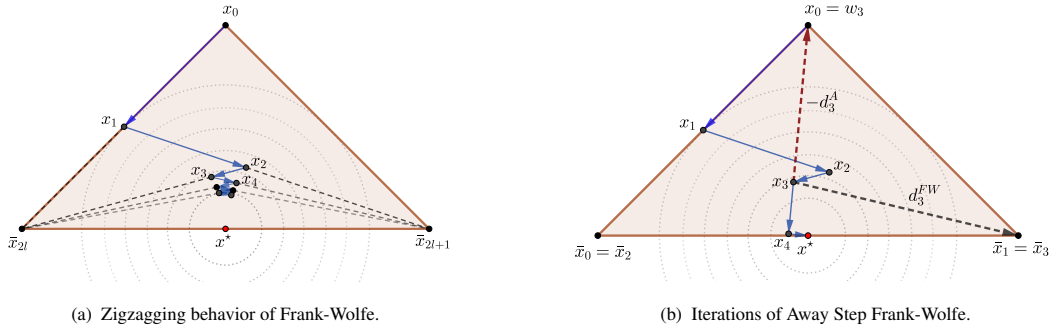


Figure 6: Example of FW and ASFW iterations.

where $\hat{d}_k := d_k / \|d_k\|$. Thus, from the last inequality, we obtain

$$h(x_k) - h(x_{k+1}) \geq \frac{(\nabla f(x_k)^T \hat{d}_k)^2}{2L}.$$

This result, along with (15), yields

$$h(x_{k+1}) \leq \left[1 - \frac{\mu}{L} \left(\frac{\nabla f(x_k)^T \hat{d}_k}{\nabla f(x_k)^T \hat{e}_k} \right)^2 \right] h(x_k). \quad (16)$$

From (16), it is clear that the linear convergence will depend on our ability to bound the ratio between $\nabla f(x_k)^T \hat{d}_k$ and $\nabla f(x_k)^T \hat{e}_k$. This can be done in certain situations for the classic FW or variants of FW using different search directions d_k . For instance, Guélat & Marcotte (1986, Theorem 2) show that if x^* is in the relative interior of C , with distance to the boundary of C at least $\rho > 0$ then, for k sufficiently large, $\nabla f(x_k)^T d_k \leq -(\rho/2) \|\nabla f(x_k)\|$ and $\gamma_k^* \leq 1$, which implies that

$$h(x_{k+1}) \leq \left[1 - \frac{\mu \rho^2}{LD^2} \right] h(x_k). \quad (17)$$

However, for the classic Frank-Wolfe, if x^* lies on the boundary of C , then $\nabla f(x_k)^T \hat{d}_k$ can get arbitrarily close to zero and the method may exhibit a zigzagging behavior (see Figure 6a).

In next section we discuss a Frank-Wolfe variant, called Away-Step Frank-Wolfe (ASFW), which exhibits linear convergence independently of the location of $x^* \in C$.

3.2. Away-Step Frank-Wolfe

We proceed to briefly reviewing ASFW algorithm presented in Lacoste-Julien & Jaggi (2015, Algorithm 1), which is summarized in Algorithm 3, and we provide a discussion on its main properties. Let V be the set of extreme points of $C := \text{conv}(\mathcal{A})$ and $n = |\mathcal{A}|$. In each iteration k of Algorithm 3, the iteration point x_k is a convex combination of elements of \mathcal{A} : $x_k = \sum_{v \in \mathcal{A}} \alpha_v^k v$. Let $\alpha^k \in \Delta_n$ be a vector containing the coefficients of such combination and define $U^k := \{v \in \mathcal{A} \mid \alpha_v^k > 0\}$. Without loss of generality, let us suppose that the elements of \mathcal{A} are ordered in such a way that x_0 is the first one: then $\alpha^0 = e_1$, the first canonical vector of \mathbb{R}^n .

Remark 3.3. In each iteration of ASFW, two search directions are computed: the FW direction d_k^{FW} and an “away direction” d_k^A . If $\nabla f(x_k)^T d_k^{FW} \leq \nabla f(x_k)^T d_k^A$, the classic FW direction is taken and the maximum allowed step-size is $\gamma_{\max} = 1$. Otherwise, $\gamma_{\max} = \alpha_{w_k} / (1 - \alpha_{w_k})$ gives a conservative upper bound for the step-size in the direction d_k^A in order to ensure $x_{k+1} \in C = \text{conv}(\mathcal{A})$; see Lacoste-Julien & Jaggi (2015) for details.

Remark 3.4. The purpose of away-steps is to reduce the weight of elements in the current active set U^k that contribute to increase the objective function (see Figure 6b). When $d_k = d_k^A$ and $\gamma_k = \gamma_{\max}$, then $U^{k+1} = U^k \setminus \{w_k\}$. If the step-size in the away direction is smaller than γ_{\max} , U_k remains unchanged, but the weight of w_k is reduced. When $d_k = d_k^{FW}$, if $\gamma_k = 1$, then $U^{k+1} = \{\bar{x}_k\}$ and $U^{k+1} = U^k \cup \{\bar{x}_k\}$ otherwise, where \bar{x}_k is given in Step 2 of Algorithm 3. This means that the cardinality of U^k either decreases, or increases by one in each iteration. Therefore, if s^k is the number of iterations, until iteration k , in which $d_k = d_k^A$ and $\gamma_k = \gamma_{\max}$, and ℓ^k is the number of iterations in which an element was added to U^k , we have $s^k + \ell^k \leq k - 1$, $s^k \leq \ell^k$ and thus $s^k \leq (k - 1)/2$.

Algorithm 3: Away-Step Frank-Wolfe (ASFW)

Data: $x_0 \in \mathcal{A}, \epsilon > 0, U^0 = \{x_0\}, \alpha^0 = e_1 \in \Delta_n$
1 for $k = 0, 1, 2, \dots$ **do**
 2 Set $\bar{x}_k \in \arg \min_{x \in \mathcal{A}} \nabla f(x_k)^T (x - x_k)$
 3 $d_k^{FW} \leftarrow \bar{x}_k - x_k$
 4 Set $w_k \in \arg \max_{x \in U^k} \nabla f(x_k)^T (x - x_k)$
 5 $d_k^A \leftarrow x_k - w_k$
 6 **if** $\nabla f(x_k)^T d_k^{FW} \geq -\epsilon$ **then** stop and return x_k
 7 **if** $\nabla f(x_k)^T d_k^{FW} \leq \nabla f(x_k)^T d_k^A$ **then**
 8 | $d_k \leftarrow d_k^{FW}, \gamma_{\max} \leftarrow 1$
 9 **else**
 10 | $d_k \leftarrow d_k^A, \gamma_{\max} \leftarrow \alpha_{w_k} / (1 - \alpha_{w_k})$
 11 **end**
 12 Set $\gamma_k \in \arg \min_{\gamma \in [0, \gamma_{\max}]} f(x_k + \gamma d_k)$
 13 $x_{k+1} \leftarrow x_k + \gamma_k d_k$.
 14 Update α^{k+1} accordingly.
 15 $U^{k+1} \leftarrow \{v \in \mathcal{A} \mid \alpha_v^{k+1} > 0\}$.
16 end

In the sequel, we present the global linear convergence of Algorithm 3. Let us define $d_k^P = \bar{x}_k - w_k$, which is known in the literature as *pairwise direction*. Clearly,

$$-\nabla f(x_k)^T d_k^P \leq -2\nabla f(x_k)^T d_k.$$

Moreover, for d_k^P we can establish a bound on the ratio $(\nabla f(x_k)^T \hat{d}_k^P) / (\nabla f(x_k)^T \hat{e}_k)$. For that, we consider the next lemma, which holds if the convex set C is a compact polyhedron, that is, $C := \{x \in \mathbb{R}^m \mid a_i^T x \leq b_i, i = 1, \dots, \ell\}$. In the following, for $U \subset \mathbb{R}^m$, define $\mathcal{I}_C(U)$ the set of active constraints of the polyhedron C for points in U , i.e., $\mathcal{I}_C(U) := \{i \in \{1, \dots, \ell\} \mid a_i^T u = b_i, \forall u \in U\}$.

Lemma 3.5 (Beck & Shtern (2017, Lemma 3.1)). *Let $C := \{x \in \mathbb{R}^m \mid a_i^T x \leq b_i, i = 1, \dots, \ell\}$ be a compact polyhedron with V being the set of its extreme points, $c \in \mathbb{R}^m$, and consider $U \subset V$. If there exists $z \in \mathbb{R}^m$ such that $a_i^T z \leq 0$, for all $i \in \mathcal{I}_C(U)$, and $c^T z > 0$, then*

$$\max_{p \in V, u \in U} c^T (p - u) \geq \frac{\Omega_C}{m+1} \frac{c^T z}{\|z\|},$$

where

$$\Omega_C := \min_{v \in V, i \in \{1, \dots, \ell\} \mid b_i > a_i^T v} \frac{b_i - a_i^T v}{\|a_i\|}.$$

Lemma 3.5 holds for the case where $C := \text{conv}(\mathcal{A})$, as $\text{conv}(\mathcal{A})$ is a polyhedron, even though we usually do not know explicitly the linear inequalities that define $\text{conv}(\mathcal{A})$. Thus, if we assume that $\bar{x}_k \in V \subset \mathcal{A}$ in Algorithm 3, we may write

$$-\nabla f(x_k)^T d_k^P = -\nabla f(x_k)^T (\bar{x}_k - w_k) = \max_{p \in V, u \in U^k} -\nabla f(x_k)^T (p - u),$$

and for $z := e_k = x^* - x_k$ and $c := -\nabla f(x_k)$, we can apply Lemma 3.5 to obtain

$$-2\nabla f(x_k)^T d_k \geq -\nabla f(x_k)^T d_k^P \geq \frac{\Omega_C}{m+1} (-\nabla f(x_k))^T \hat{e}_k > 0,$$

which shows that, in each iteration of ASFW for which $\gamma_k^* \leq \gamma_{\max}$, we have

$$h(x_{k+1}) \leq \left[1 - \frac{\mu}{4L} \left(\frac{\Omega_C}{D(m+1)} \right)^2 \right] h(x_k).$$

When $d_k = d_k^A$ and $\gamma_k^* > \gamma_{\max}$, we can only ensure that $h(x_{k+1}) \leq h(x_k)$. However, as discussed before (see Remark 3.4), this type of iteration cannot occur too often. Therefore, for the ASFW we have

$$h(x_k) \leq h(x_0) \left[1 - \frac{\mu}{4L} \left(\frac{\Omega_C}{D(m+1)} \right)^2 \right]^{(k-1)/2}, \quad \forall k \in \mathbb{N}, \quad (18)$$

which implies global linear convergence (since $\Omega_C \leq D$ and $\mu \leq L$). We summarize the above discussion, which combines the analysis of Lacoste-Julien & Jaggi (2015) with Lemma 3.5, in the following result.

Theorem 3.6. *Let f be μ -strongly convex, ∇f L -Lipschitz, C a compact polyhedral set (with diameter D) and let f^* be the optimal value for (13). Let $(x_k)_{k \in \mathbb{N}}$ be the sequence generated by ASFW for problem (13) and assume that $\{\bar{x}_k\} \subset V$. Then, for all $k \geq 1$*

$$f(x_k) - f^* \leq (1 - \eta)^{(k-1)/2} (f(x_0) - f^*),$$

where $\eta := \frac{\mu}{4L} \left(\frac{\Omega_C}{D(m+1)} \right)^2$, with Ω_C from Lemma 3.5.

Remark 3.7. Note that, for problem (4), $f \equiv \Psi$, and such function satisfies the assumptions in Theorem 3.6, with $\mu = L = 1$.

3.3. Frank-Wolfe similarities with the Triangle Algorithm

If we apply the Frank-Wolfe algorithm (Algorithm 2) to formulation (4), in each iteration, we must solve the subproblem

$$\begin{aligned} \min_y \quad & (y_k - p)^T (y - y_k) \\ \text{s.t.} \quad & y \in \mathcal{A}. \end{aligned} \quad (19)$$

Since $\mathcal{A} = \{v_1, \dots, v_n\}$ is a finite set, \bar{y}_k solution of (19) can be chosen as $\bar{y}_k = v_j$, where $v_j^T (y_k - p) = \min_{1 \leq i \leq n} v_i^T (y_k - p)$. Taking into account that $y_k \in \text{conv}(\mathcal{A})$ for every k , we may associate y_k with iterate p_k of Algorithm 1 to obtain

$$v_j^T (p_k - p) = \min_{1 \leq i \leq n} v_i^T (p_k - p). \quad (20)$$

In comparison, in the k -th iteration, Algorithm 1 chooses $v_i \in \mathcal{A}$ such that $v_i^T (p_k - p) \leq (\|p_k\|^2 - \|p\|^2)/2$, which, by Lemma 2.2, characterizes a pivot. Thus, it is clear by (20), that Algorithm 2 not only chooses a pivot, but one that minimizes the product $v_i^T (p_k - p)$. In this sense, Frank-Wolfe (Algorithm 2) applied to (4), with *exact line-search*, can be seen as a *greedy* version of the Triangle Algorithm.

On the other hand, the Triangle Algorithm can be interpreted as an *inexact* version of Frank-Wolfe, in the following sense. Similar to Jaggi (2013), we consider $\bar{y}_k \in C$ an inexact solution of the FW subproblem (see Step 2 in Algorithm 2, replacing x by y , and f by Ψ), if the following inequality holds

$$(\bar{y}_k - y_k)^T \nabla \Psi(y_k) \leq \min_{y \in C} (y - y_k)^T \nabla \Psi(y_k) + \frac{1}{2} \hat{\delta}_k \bar{\gamma}_k C_\Psi,$$

where $\hat{\delta}_k \geq 0$ is an inexactness parameter, $\bar{\gamma}_k$ is given in (6) and C_Ψ is the *curvature constant* of Ψ with respect to the set C (see Jaggi (2013) for a formal definition and detailed discussion). If we consider (4), we have $\Psi(y) = \|y - p\|^2/2$, $\nabla \Psi(y) = y - p$, and $C = \text{conv}(\mathcal{A})$ (for which $C_\Psi = D^2$). Then, by identifying y_k with p_k and setting the inexactness parameter as

$$\hat{\delta}_k = \frac{2}{D^2} \|v_j - p_k\|^2, \quad (21)$$

where $j \in \arg \min\{v_i^T (p_k - p) \mid 1 \leq i \leq n\}$, if v_i is a pivot, we have

$$\begin{aligned} \min_{y \in C} (y - p_k)^T \nabla \Psi(p_k) + \frac{1}{2} \bar{\gamma}_k C_\Psi \hat{\delta}_k &= (v_j - p_k)^T (p_k - p) + \frac{1}{2} \bar{\gamma}_k D^2 \frac{2}{D^2} \|v_j - p_k\|^2 \\ &= (v_j - p_k)^T (p_k - p) - \frac{(p_k - p)^T (v_j - p_k)}{\|v_j - p_k\|^2} \|v_j - p_k\|^2 = 0 \\ &\geq -\frac{1}{2} \|p_k - p\|^2 \geq (p_k - p)^T (v_i - p_k), \end{aligned} \quad (22)$$

Table 1: Iterations complexity to find $p_k \in \text{conv}(\mathcal{A})$ such that $d(p_k, p) < \varepsilon R$. The last assumption applies only to TA.

Assumptions	Triangle Algorithm	Frank-Wolfe	ASFW
$p \in \text{conv}(\mathcal{A})$	$\mathcal{O}(1/\varepsilon^2)$	$\mathcal{O}(D^2/\varepsilon^2)$	$\mathcal{O}\left(\frac{D^2(m+1)^2}{\Omega_C^2} \ln \frac{\delta_0}{\varepsilon}\right)$
and Assumption 2.13	$\mathcal{O}(1/\varepsilon^2)$	$\mathcal{O}\left(\frac{D^2}{\rho^2} \ln \frac{\delta_0}{\varepsilon}\right)$	$\mathcal{O}\left(\frac{D^2(m+1)^2}{\Omega_C^2} \ln \frac{\delta_0}{\varepsilon}\right)$
and Assumption 2.14	$\mathcal{O}\left(\frac{R^2}{\rho^2} \ln \frac{\delta_0}{\varepsilon R}\right)$	$\mathcal{O}\left(\frac{D^2}{\rho^2} \ln \frac{\delta_0}{\varepsilon}\right)$	$\mathcal{O}\left(\frac{D^2(m+1)^2}{\Omega_C^2} \ln \frac{\delta_0}{\varepsilon}\right)$

where in the second inequality we use Lemma 2.2(iii). Thus, (22) shows that pivot v_i in the k -th iteration of Algorithm 1, is an inexact solution to the FW subproblem for $\hat{\delta}_k$ as in (21).

Therefore, it is not surprising that the complexity results for both algorithms are alike. As discussed in Section 3 (see inequality (14)), the iterations of the Frank-Wolfe method for the problem (13) satisfy

$$\Psi(y_k) - \Psi(y^*) = \mathcal{O}\left(\frac{1}{k}\right).$$

In particular, for the convex hull membership problem with $p \in \text{conv}(\mathcal{A})$, it means

$$\frac{1}{2}\|y_k - p\|^2 = \mathcal{O}\left(\frac{1}{k}\right).$$

Furthermore, by Theorem 2.6, when $p \in \text{conv}(\mathcal{A})$, the Triangle Algorithm generates a sequence of points $p_k \in \text{conv}(\mathcal{A})$ such that

$$\|p_k - p\| < \sqrt{\frac{48}{k}}R,$$

thereby exhibiting the same iteration complexity as Frank-Wolfe, if we identify $y_k \equiv p_k$.

Interestingly, the iteration complexity is improved for both algorithms when p is in the relative interior of $\text{conv}(\mathcal{A})$ as stated in Theorem 2.16 for the Triangle algorithm with strict pivots and as can be devised from (17) for Frank-Wolfe (see also Guélat & Marcotte (1986, Theorem 2)). In Table 1, we present a comparison between the iteration complexity of the Triangle Algorithm for the case when $p \in \text{conv}(\mathcal{A})$ with those of Frank-Wolfe and Away Step Frank-Wolfe (ASFW) under different assumptions. The constants in the table for FW and ASFW are derived from (14), (16) and (17).

In addition, when $p \notin \text{conv}(\mathcal{A})$, according to Theorem 2.6(ii), the iteration complexity of the Triangle Algorithm depends only on the ‘‘problem geometry’’, that is, it does not depend on the tolerance ε . Since FW can be viewed as a Greedy TA, this assertion also holds, as long as stopping criteria based on distance duality are integrated to FW (see Section 5).

When it comes to the cost per iteration, although in the worst case both algorithms have to evaluate the entire list of inner products $v_i^T(p_k - p)$, it is expected that, in general, the iteration of the Algorithm 1 requires fewer evaluations because it only asks for a simple pivot, in contrast to the Frank-Wolfe method that needs a pivot which minimizes the inner product $v_i^T(p_k - p)$.

3.4. Frank-Wolfe as a Greedy Triangle algorithm

We have just investigated the relationship between the Triangle Algorithm and the Frank-Wolfe method applied to the problem (4): FW is equivalent to the Triangle Algorithm when in TA one chooses a pivot v_j such that

$$v_j \in \arg \min\{v_i^T(p_k - p) \mid v_i \in \mathcal{A} \setminus \{p_k\}\}. \quad (23)$$

In view of Lemma 2.2(ii), problem (23) means that the Frank-Wolfe method (with the peculiarities of the CHMP) is equivalent to a *greedy* Triangle Algorithm. However, instead of choosing v_j such that

$$v_j^T(p_k - p) \leq \frac{\|p_k\|^2 - \|p\|^2}{2}, \quad (24)$$

we select $v_j \in \mathcal{A} \setminus \{p_k\}$ such that the left side of (24) is minimal. We call this algorithm *Greedy Triangle* (GT) and summarize it in Algorithm 4.

Algorithm 4: GREEDY TRIANGLE (GT)

Data: $\mathcal{A}, p, R, \varepsilon \in (0, 1)$
1 Choose $p_0 \in \arg \min\{d(v_j, p) \mid v_j \in \mathcal{A}\}$
2 **for** $k = 0, 1, 2, \dots$ **do**
3 **if** $d(p_k, p) < \varepsilon R$ **then** stop.
4 Choose $v_j \in \arg \min\{v_i^T(p_k - p) \mid v_i \in \mathcal{A} \setminus \{p_k\}\}$.
5 **if** $v_j^T(p_k - p) > p^T(p_k - p)$ **then** stop: p_k is a p -witness.
6 Set $\tilde{\gamma}_k \in \arg \min\{d(p, (1 - \gamma)p_k + \gamma v_j) \mid 0 \leq \gamma \leq 1\}$.
7 $p_{k+1} \leftarrow (1 - \tilde{\gamma}_k)p_k + \tilde{\gamma}_k v_j$.
8 **end**

Remark 3.8. We point out that, for $p = 0$, Algorithm 4 coincides with the von Neumann² algorithm (Epelman & Freund, 2000; Gonçalves et al., 2009; Peña et al., 2016).

Since the inequality $v_j^T(p_k - p) \leq p^T(p_k - p)$ characterizes a strict pivot, in view of Step 4 of Algorithm 4 it follows that a strict pivot is chosen whenever possible. Furthermore, if there are strict pivots at p_k verifying condition (12) (see Assumption 2.14), Algorithm 4 will employ one of them.

Proposition 3.9. *Assume that Assumption 2.13 holds.*

- (i) *If $d(p_k, p) > \varepsilon R$, then there exists a strict pivot $w \in \mathcal{A}$ satisfying (12).*
- (ii) *Algorithm 4 always chooses $v_j \in \mathcal{A}$ satisfying condition (12).*

Proof. (i) Let $p_k \in \text{conv}(\mathcal{A})$ be such that $d(p, p_k) > \varepsilon R$. Assumption 2.13 implies that there exists $\tau > 0$ such that $q = p + \tau(p - p_k) \in \text{conv}(\mathcal{A})$ and $d(q, p) = \rho$. Since $q \in \text{conv}(\mathcal{A})$, applying Theorem 2.9 for q (instead of p), we conclude that there exists a strict q -pivot w at p . Thus, inequality (12) follows from Remark 2.8.

(ii) Consider v_j of Step 4 in Algorithm 4 and let $w \in \mathcal{A}$ be a strict pivot satisfying (12). Then,

$$\begin{aligned} (p - q)^T(v_j - q) &= -\tau(p - p_k)^T(v_j - p - \tau(p - p_k)) = -\tau(p - p_k)^T v_j + \tau(p - p_k)^T p + \tau^2 \|p - p_k\|^2 \\ &\leq -\tau(p - p_k)^T w + \tau(p - p_k)^T p + \tau^2 \|p - p_k\|^2 = (p - q)^T(w - q) \leq 0. \end{aligned}$$

Hence, v_j also satisfies the inequality (12). □

A possible drawback of Algorithm 4, in comparison with Algorithm 1, is that, at each iteration, it is necessary to go through the entire list \mathcal{A} to assure the minimum of (23). Therefore, the cost per iteration is always $O(nm)$.

4. Projected gradient methods for CHMP

Another class of first-order methods that can be used to solve the CHMP is that of projected gradient methods. The projected gradient method is well known in the literature to solve problems as (13). Given a feasible point $x_k \in C$, the method considers the search direction $d_k := \bar{x}_k - x_k$, where $\bar{x}_k := P_C(x_k - \nabla f(x_k))$ and P_C denotes the orthogonal projection onto the closed and convex set C . It is not difficult to show that if x_k is not a stationary point for (13), then d_k is a descent direction, and that $x_k + \gamma d_k \in C$, for all $\gamma \in [0, 1]$. On the other hand, one can show that $\|d_k\| = 0$ if, and only if, x_k is a stationary point for (13).

Over the last decades, nonmonotone strategies for nonlinear optimization began to become popular (Grapiglia & Sachs, 2017; Grippo et al., 1986; Zhang & Hager, 2004) and, along with scaling ideas, are present in enhanced versions of the projected gradient method. An example is the Spectral Projected Gradient (SPG) method, proposed by Birgin et al. (2000, 2009), which uses Barzilai-Borwein (spectral, BB-type) scaling of the gradient direction combined with nonmonotone line-search.

²According to Dantzig, von Neumann communicated that algorithm to him in the 1940s and Dantzig studied it later in an unpublished manuscript (Dantzig, 1992).

The SPG iteration is given by $x_{k+1} = x_k + \gamma_k d_k$, where the search direction d_k is defined as $d_k := P_C(x_k - \lambda_k \nabla f(x_k)) - x_k$. Here, λ_k is the spectral scaling parameter

$$\lambda_k = \frac{s_{k-1}^T s_{k-1}}{s_{k-1}^T u_{k-1}}, \quad (25)$$

where $s_{k-1} := x_k - x_{k-1}$ and $u_{k-1} := \nabla f(x_k) - \nabla f(x_{k-1})$. Rather than imposing a sufficient decrease at each iteration, a characteristic of the SPG is to employ a nonmonotone line-search in order to favor the acceptance of the full-step $\gamma_k = 1$, while still guaranteeing global convergence (Birgin et al., 2000, Section 2).

For instance, the nonmonotone line-search proposed by Grippo et al. (1986) depends on an integer parameter $M \geq 1$ and imposes a functional decrease with respect to the highest functional value over the last M iterations (if $M = 1$ the line-search is monotone).

Algorithm 5 summarizes the Spectral Projected Gradient method as in Birgin et al. (2014).

Algorithm 5: SPECTRAL PROJECTED GRADIENT (SPG)

Data: $x_0 \in C, \epsilon > 0, M \geq 1, \eta \in (0, 1), 0 < \lambda_{\min} \leq \lambda_{\max} < \infty$ and $\lambda_0 \in [\lambda_{\min}, \lambda_{\max}]$.

```

1 for  $k = 0, 1, 2, \dots$  do
2    $\bar{x}_k = P_C(x_k - \lambda_k \nabla f(x_k))$  and define  $d_k = \bar{x}_k - x_k$ .
3   if  $\|d_k\| < \epsilon$  then stop: return  $x_k$ .
4   Compute  $f_{\max} = \max\{f(x_{k-j}) \mid 0 \leq j \leq \min\{k, M-1\}\}$  and set  $\gamma \leftarrow 1$ .
5   while  $f(x_k + \gamma d_k) > f_{\max} + \eta \gamma \nabla f(x_k)^T d_k$  do
6      $\gamma \leftarrow \gamma/2$ 
7   end
8    $\gamma_k \leftarrow \gamma$ 
9    $x_{k+1} \leftarrow x_k + \gamma_k d_k$ .
10   $s_k \leftarrow x_{k+1} - x_k$ 
11   $u_k \leftarrow \nabla f(x_{k+1}) - \nabla f(x_k)$ .
12  if  $s_k^T u_k \leq 0$  then
13     $\lambda_{k+1} \leftarrow \lambda_{\max}$ .
14  else
15     $\lambda_{k+1} \leftarrow \max\{\lambda_{\min}, \min\{s_k^T s_k / s_k^T u_k, \lambda_{\max}\}\}$ .
16  end
17 end

```

Remark 4.1. The main cost per iteration of projected gradient methods is the cost of computing the projection $P_C(\cdot)$. Fortunately, if we consider formulation (3) for the CHMP, then $f \equiv \Phi$ and $C = \Delta_n$, the unit simplex, for which the projection can be computed in $O(n \log n)$ operations (Condat, 2016).

5. Specific stopping criteria for CHMP

The distance duality (see Theorem 1.1) is a remarkable result for the CHMP, and it is the keystone of the Triangle Algorithm. Thus, a natural question that arises is: how other first-order methods applied to CHMP can benefit from this result in order to employ specialized stopping criteria?

From the discussion in Section 3.4, we established that Frank-Wolfe applied to CHMP can be seen as a Greedy Triangle Algorithm, in the sense that it chooses $v_j \in \mathcal{A} \setminus \{p_k\}$ that minimizes $v_j^T (p_k - p)$. Hence, in case $v_j^T (p_k - p) > (\|p_k\|^2 - \|p\|^2)/2$, we can assert that there is no pivot at p_k and $p \notin \text{conv}(\mathcal{A})$. This way, as in the Triangle Algorithm, it is possible to stop FW (and its variants) as soon as a witness is detected. Furthermore, as in TA, we also stop FW iterations as soon as $\|p_k - p\| \leq \epsilon R$.

Another important question is how to set the tolerance ϵ in Algorithm 2 (Step 4) in order to ensure $\|p_k - p\| \leq \epsilon R$. Note that the classic stopping criterion of Algorithm 2 (and Algorithm 3) is $\nabla \Psi(y_k)^T (\bar{y}_k - y_k) \geq -\epsilon$ (where we denote an iterate by y_k instead of x_k because we are considering problem (4)). Let $y^* \in \text{conv}(\mathcal{A})$ be the solution of (4). From the discussion in Section 3, we know that

$$\Psi(y_k) - \Psi(y^*) \leq \nabla \Psi(y_k)^T (y_k - \bar{y}_k) \leq \epsilon,$$

where \bar{y}_k is the solution of the FW subproblem. For problem (4), this inequality is equivalent to

$$\|y_k - p\|^2 - \|y^* - p\|^2 \leq 2\epsilon,$$

which yields $\|y_k - p\| - \|y^* - p\| \leq 2\epsilon/\|y_k - p\|$. Thus, imposing $\epsilon := \|y_k - p\|\epsilon R/2$, we have

$$\|y_k - p\| \leq \|y^* - p\| + \epsilon R. \quad (26)$$

If $p \in \text{conv}(\mathcal{A})$, y^* , solution of (4), is such that $\|y^* - p\| = 0$ and then inequality (26) coincides with $\|p_k - p\| \leq \epsilon R$ (with $p_k = y_k$). On the other hand, if $\nabla\Psi(y_k)^T(y_k - \bar{y}_k) \leq \|y_k - p\|\epsilon R/2$ and $\|y_k - p\| > \epsilon R$ then, from (26), $\|y^* - p\| > 0$ and we can conclude that $p \notin \text{conv}(\mathcal{A})$. Therefore, instead of the classic FW stopping criterion, we can use the *relative error* criterion

$$\nabla\Psi(y_k)^T(y_k - \bar{y}_k) \leq \|y_k - p\| \frac{\epsilon R}{2}.$$

When it comes to projected gradient methods applied to problem (3), we recall that $\nabla\Phi(x_k) = A^T(Ax_k - p)$ and, since $x_k \in \Delta_n$, we can write $p_k = Ax_k \in \text{conv}(\mathcal{A})$. Thus, each gradient component can be written as $[\nabla\Phi(x_k)]_i = v_i^T(p_k - p)$ and, in case each component is greater than $(\|p_k\|^2 - \|p\|^2)/2$, we can stop and return $p_k = Ax_k$ as a witness that $p \notin \text{conv}(\mathcal{A})$.

Next, we discuss how to set the tolerance ϵ in Algorithm 5 (Step 3) in order to guarantee that $\|A\bar{x}_k - p\| \leq \epsilon R$ (here, \bar{x}_k is from Step 2). For that, first we need the following lemma.

Lemma 5.1. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuously differentiable convex function with L -Lipschitz gradient and $C \subset \mathbb{R}^n$ be a non-empty, convex and compact set with diameter D . If $x^* \in C$ is a minimizer of f over C and $\bar{x} := P_C(x - \lambda\nabla f(x))$, for some $\lambda > 0$ and $x \in C$, then*

$$f(\bar{x}) - f(x^*) \leq D \left(\frac{L}{2} + \frac{1}{\lambda} \right) \|\bar{x} - x\|. \quad (27)$$

Proof. Since \bar{x} is the projection of $x - \lambda\nabla f(x)$ onto C , by the characterization of projections onto convex sets we have $(w - \bar{x})^T(x - \lambda\nabla f(x) - \bar{x}) \leq 0$, for all $w \in C$. In particular, as $x^* \in C$ and $\lambda > 0$, we obtain

$$\nabla f(x)^T(\bar{x} - x^*) \leq \frac{1}{\lambda}(\bar{x} - x^*)^T(x - \bar{x}). \quad (28)$$

Note also that

$$(\bar{x} - x^*)^T(x - \bar{x}) = (\bar{x} - x)^T(x - \bar{x}) + (x - x^*)^T(x - \bar{x}) \leq -\|x - \bar{x}\|^2 + \|x - x^*\| \|x - \bar{x}\|. \quad (29)$$

Therefore, using (28) and (29), we get, after some manipulations,

$$\nabla f(x)^T(\bar{x} - x) \leq \nabla f(x)^T(x^* - x) + \frac{1}{\lambda}(-\|x - \bar{x}\|^2 + \|x - x^*\| \|x - \bar{x}\|). \quad (30)$$

Since ∇f is L -Lipschitz, by using (30) and convexity of f , we obtain

$$f(\bar{x}) - f(x) - \frac{L}{2} \|\bar{x} - x\|^2 \leq f(x^*) - f(x) + \frac{1}{\lambda}(-\|x - \bar{x}\|^2 + \|x - x^*\| \|x - \bar{x}\|),$$

which implies that

$$f(\bar{x}) - f(x^*) \leq \frac{L}{2} \|\bar{x} - x\|^2 + \frac{1}{\lambda}(-\|x - \bar{x}\|^2 + \|x - x^*\| \|x - \bar{x}\|) \leq D \left(\frac{L}{2} + \frac{1}{\lambda} \right) \|\bar{x} - x\|,$$

where, in the last inequality, we used the fact that both $x, x^* \in C$, which has, by hypothesis, diameter D . \square

Now, we recall that for problem (3), $f(x) := \Phi(x) = \frac{1}{2} \|Ax - p\|^2$ and $C = \Delta_n$. In this case, $L = \|A\|^2$ and $D = \sqrt{2}$. From the definition of the spectral parameter (see (25)), we have $1/\lambda_k \leq L$. Then, we can rewrite inequality (27) from Lemma 5.1 as

$$\Phi(\bar{x}_k) - \Phi(x^*) \leq \frac{3}{2} LD \|\bar{x}_k - x_k\|, \quad (31)$$

by taking x_k as x and \bar{x}_k as \bar{x} .

Since, in Algorithm 5, $d_k := \bar{x}_k - x_k$, if one imposes $\epsilon := \frac{\|A\bar{x}_k - p\|}{3LD} \epsilon R$, inequality (31) implies that $\|A\bar{x}_k - p\| \leq \|Ax^* - p\| + \epsilon R$. If $p \in \text{conv}(\mathcal{A})$, x^* , solution of (3), is such that $\|Ax^* - p\| = 0$ and then the last inequality coincides with $\|A\bar{x}_k - p\| \leq \epsilon R$. We remark that computing a \bar{x}_k such that the last inequality is satisfied is similar to finding an ϵ -solution for the Triangle Algorithm because $A\bar{x}_k \in \text{conv}(\mathcal{A})$. On the other hand, if $\|d_k\| \leq \frac{\|A\bar{x}_k - p\|}{3LD} \epsilon R$ and $\|A\bar{x}_k - p\| > \epsilon R$, then $\|Ax^* - p\| > 0$ and we can conclude that $p \notin \text{conv}(\mathcal{A})$.

6. Numerical Experiments

In order to evaluate the performance of the first-order methods described in Sections 2, 3 and 4 when solving CHMP, equipped with the stopping criteria from Section 5, we present three sets of computational experiments.

In Section 6.1, we consider randomly generated instances of CHMP organized in 4 different scenarios according to the distribution of the points of \mathcal{A} and the relative position of the query point p . This way, we can cover different ranges for values of the *geometric constants*, namely $D, R, \rho, \delta_0, \Delta$ and c , and evaluate the performance of the algorithms in favorable and unfavorable conditions. Section 6.2 discuss how linear programming feasibility problems can be cast as CHMPs under suitable assumptions and compares the studied first-order methods for CHMP with the dual-simplex applied to the original problem. Finally, in Section 6.3 we apply the studied methods to an image classification problem, in order to determine membership of (or distance from) elements in the testing set to the convex hull of the elements in the training set (or its subclasses).

All the tests were run on an Intel Core i7-8565U 1.80 GHz with 8 GB of RAM running Windows 10 whilst the algorithms were implemented in Matlab R2019a. The codes of our experiments are fully available at <https://github.com/rafaelafilipozzi/First-order-methods-for-the-CHMP>.

6.1. Artificial instances of CHMP

The artificial instances of CHMP were generated following a procedure similar to the one proposed by Li & Kalantari (2013). Each element of $\mathcal{A} = \{v_1, \dots, v_n\}$ is randomly generated according to a uniform distribution in the unit ball of \mathbb{R}^m (Harman & Lacko, 2010). For that, each of the n points of \mathcal{A} is generated as follows: first we sample a $\hat{v} \in \mathbb{R}^m$ from the standard normal distribution; then we set $v \in \mathcal{A}$ as $v := \sqrt{u}(\hat{v}/\|\hat{v}\|)$, where u is sampled from the uniform distribution in $[0, 1]$.

The query point $p \in \mathbb{R}^m$ is generated in the following four cases:

- (a) $p \in \text{conv}(\mathcal{A})$ in the relative interior of $\text{conv}(\mathcal{A})$;
- (b) $p \in \text{conv}(\mathcal{A})$ with visibility factor close to zero;
- (c) $p \notin \text{conv}(\mathcal{A})$ away from the boundary;
- (d) $p \notin \text{conv}(\mathcal{A})$ with visibility factor close to zero.

The cases above were selected aiming to evaluate the performance of the considered first-order algorithms in the best and worst cases for the Triangle Algorithm. For case (a) we point out the improved iteration complexity given by Theorem 2.16 for the Triangle Algorithm with “strong pivots” (see Assumptions 2.13, 2.14 and Proposition 3.9). In the cases where the point p is outside $\text{conv}(\mathcal{A})$, the relative position of p affects the values of Δ, R and δ_0 in Theorems 2.6 and 2.11. In the cases (b) and (d) we have tried to force the visibility factor c to be close to zero in order to obtain harder instances for the Triangle algorithms (TA and GT). Details are given in Sections 6.1.2 and 6.1.4.

For all the tests in this section we consider the following tolerances, algorithmic parameters and implementation details:

- (i) The maximum number of iterations was set to $\text{maxit} := \min\{\max\{1000n, 10000\}, 10^6\}$.
- (ii) The usual stopping criteria of both Away Step Frank-Wolfe (ASFW) and Spectral Projected Gradient (SPG) have been changed according to the discussions in Section 5.
- (iii) The tolerance for an ϵ -approximate solution of CHMP was set to $\epsilon := 10^{-4}$ for all tested algorithms.

- (iv) For the nonmonotone line search in SPG we used $M = 15$ because it obtained the best performance in preliminary tests. Moreover, we set $\lambda_{\min} := 10^{-8}$ and $\lambda_{\max} := 10^8$.
- (v) As the variables $x \in \mathbb{R}^n$ in the quadratic formulation (3) are the coefficients of the convex combination (of the elements of \mathcal{A}), we consider the starting point $x_0 = e_i \in \mathbb{R}^n$, where e_i denote the i -th canonical vector of \mathbb{R}^n and the index i is such that

$$i = \arg \min_{1 \leq j \leq n} \{\|v_j - p\|\},$$

so $Ax_0 = p_0$, where p_0 is the starting point for Algorithm 1.

- (vi) In our implementation of the Triangle Algorithm, at each iteration, we build the full list of values $v_i^T(p_k - p)$ (see Lemma 2.2(ii)), find all the pivots (if any) and select one randomly. The random choice of a pivot (instead of a pivot v_i with the smallest index i) is inspired by Zhang & Kalantari (2016) and turns out to be very effective in reducing the number of TA iterations across the four different scenarios considered in our experiments. Though we could have tested sequentially whether v_i , for $i = 1, 2, \dots, n$, is a pivot and evaluated the products $v_i^T(p_k - p)$ one-by-one, stopping the search as soon as the first pivot is found, we instead “vectorize” the search which usually has a superior performance in Matlab than loops. We emphasize that in TA we have used simple pivots (see Lemma 2.2) at each iteration rather than strict pivots as in Definition 2.7.

The computational cost per iteration of both TA and GT is $O(mn)$ arithmetic operations, corresponding to n inner products of vectors in \mathbb{R}^m . In ASFW, we also have $O(mn)$ arithmetic operations per iteration to determine Frank-Wolfe and Away directions. SPG is the algorithm that has the highest computational cost per iteration, $O(2mn + n \log n)$, since it requires two matrix-vector products and the projection onto the unit simplex Δ_n ; see Section 4.

Remark 6.1. Concerning the Triangle Algorithm, it was shown in (Kalantari, 2019a) that, with a pre-processing cost of $O(mn^2)$, the iteration cost of TA can be reduced to $O(m + n)$. Such alternative implementation may be interesting when the expected number of TA iterations is much larger than n . However, in the experiments reported here we do not use pre-processing.

In the experiments described in this subsection we fix the dimension as $m = 100$ and vary the number of points n in \mathcal{A} from 500 to 10^5 . For each value of n , 10 random instances are generated. Tables in Appendix A show the *average* number of iterations required by each algorithm.

6.1.1. Case(a): p in the relative interior of $\text{conv}(\mathcal{A})$

To ensure that the query point lies in $\text{relint}(\text{conv}(\mathcal{A}))$, we select it as $p = 0 \in \mathbb{R}^m$, the center of the unit ball. Since every $v_j \in \mathcal{A}$ originates from a uniform distribution on the unit ball, whenever the number n of points in \mathcal{A} is sufficiently large, it is rather true that $p \in \text{conv}(\mathcal{A})$. In addition, it is also likely that $B_\rho(p) \subset \text{conv}(\mathcal{A})$, for $\rho > 0$ away from zero and thus, in this case, Assumption 2.13 is satisfied.

In this scenario, we expect GT to perform better than TA. Bear in mind, from Theorems 2.6 and 2.11, that the iteration complexity of TA is given by

$$O\left(\min\left\{\frac{1}{\varepsilon^2}, \frac{1}{c} \ln \frac{\delta_0}{\varepsilon R}\right\}\right),$$

where $\delta_0 = d(p, p_0)$, $R := \max\{d(v_i, p) \mid v_i \in \mathcal{A}\}$ and c is from Assumption 2.10. However, in view of Proposition 3.9, Assumption 2.14 is satisfied by GT and its iteration complexity is given by Theorem 2.16:

$$O\left(\frac{R^2}{\rho^2} \ln \frac{\delta_0}{\varepsilon R}\right).$$

Figure 7 shows the performance of the algorithms in terms of running time, while the number of iterations is listed in Table A.5 (in the appendix). As anticipated, GT requires fewer iterations than TA to retrieve p_ε such that $d(p_\varepsilon, p) \leq \varepsilon R$. Moreover, when n is large, GT also achieves the best performance in terms of running time, followed closely by SPG and ASFW. Even though the cost per iteration of SPG is higher, it required only 12 iterations on average while GT and ASFW used up to 118 iterations.

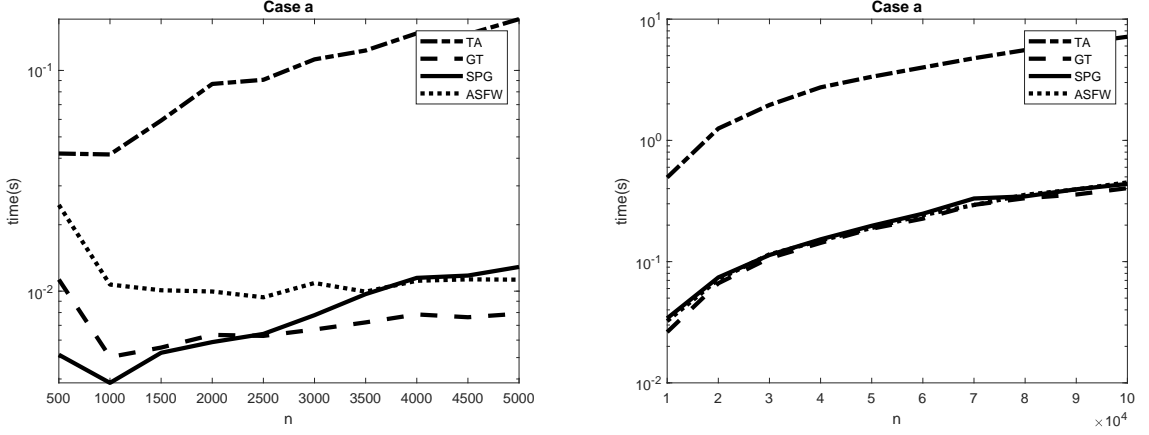


Figure 7: Case (a). Running times (in seconds) for dimension $m = 100$ and increasing n , the number of points in \mathcal{A} .

6.1.2. Case (b): $p \in \text{conv}(\mathcal{A})$ with visibility factor close to zero

Here, we generate the query point $p \in \text{conv}(\mathcal{A})$ on the boundary of the convex hull, similar to Example 2.12. For that, first we determine two points of \mathcal{A} , say v_ℓ and v_q , that have the highest values of the linear functional $\psi(v) = e^T v$: since $\text{conv}(\mathcal{A})$ is a polytope and from the way the points of \mathcal{A} were randomly generated, we expect, for a sufficiently large n , that v_ℓ and v_q are extreme points of $\text{conv}(\mathcal{A})$ (the maximum of $\psi(v)$ is achieved in at least one of these points). Next, we set $p = v_\ell/2 + v_q/2$. In addition, to prevent v_ℓ and v_q from being points of \mathcal{A} closest to p , we add to \mathcal{A} a new point v_s such that $d(v_s, p) < d(v_\ell, v_q)/2$. More precisely,

$$v_s = p - \frac{\beta \|v_\ell - v_q\|}{2 \|p\|} p,$$

where $\beta = 0.9$. This last step is necessary because if $p_0 = v_q = \text{argmin}\{d(v, p) \mid v \in \mathcal{A}\}$, for example, v_ℓ is a strict pivot since

$$(p_0 - p)^T (v_q - p) = (v_\ell - p)^T (v_q - p) = \frac{1}{4} (v_\ell - v_q)^T (v_q - v_\ell) \leq 0.$$

In this situation, GT could select v_ℓ as pivot in the first iteration and p would be retrieved after the exact line search. This actually happened when we did not add v_s .

The instances considered in this subsection are more difficult for TA and GT algorithms because the upper bound $1/\sqrt{1+c}$ in (10) can be quite close to 1. Since $p = \frac{1}{2}(v_q + v_\ell)$, either v_q or v_ℓ is a pivot for any $p' \in \text{conv}(\mathcal{A}) \setminus \{p\}$. In fact, otherwise

$$(p' - p)^T (v_\ell - p) > \frac{1}{2} \|p' - p\|^2 \quad \text{and} \quad (p' - p)^T (v_q - p) > \frac{1}{2} \|p' - p\|^2$$

would imply

$$0 = (p' - p)^T (v_\ell + v_q - 2p) > \|p' - p\|^2,$$

a contradiction. Moreover, if v_ℓ was the pivot chosen to obtain p_k at the iteration $k - 1$, from the exact line search $(v_\ell - p_k)^T (p - p_k) = 0$, thus v_ℓ is not a pivot at iteration k , which implies that v_q is. Therefore, starting at $p_0 \in \text{conv}(\mathcal{A}) \setminus \{v_\ell, v_q, p\}$, we could have in each iteration of TA (or GT) p_k as a convex combination of p_0 , v_ℓ and v_q , with v_q and v_ℓ alternating as pivots. This phenomenon is similar to the one of Example 2.12. Hence, since $1/c$ can be relatively large, we expect to observe the complexity of the Theorem 2.6 instead of the Theorem 2.11 for TA and GT (in contrast to case (a)). Nevertheless, the assumptions of Theorem 3.6, concerning ASFW, are still satisfied.

Both Triangle and Greedy Triangle algorithms reached the maximum number of iterations and, therefore, are not reported in Table A.5 and Figure 8. On average, the ASFW converged in 12 iterations and 7 of such

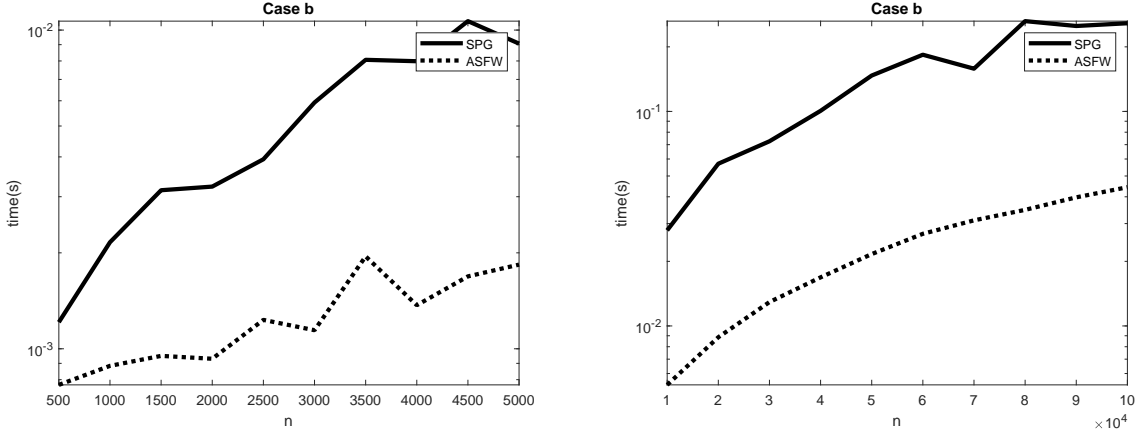


Figure 8: Case (b). Running times (in seconds) for dimension $m = 100$ and increasing n , the number of points in \mathcal{A} .

iterations have used the away direction. From Figure 8 we see that ASFW outperformed SPG. This happened because, although SPG achieved the desired precision in a maximum of 10 iterations, on average (see Table A.5), its iteration cost is higher than ASFW.

6.1.3. Case (c): $p \notin \text{conv}(\mathcal{A})$ and away from the boundary

Now, the query point p is located outside $\text{conv}(\mathcal{A})$. We generate $p \in \mathbb{R}^m$ as in Case (b), followed by a dilation: it is multiplied by 1.5. Differently from the previous case, we do not add another point v_s to $\text{conv}(\mathcal{A})$ to prevent p_0 from being v_q or v_ℓ . Note that, in this way, as $\mathcal{A} \subset B_1(0)$, $p \notin \text{conv}(\mathcal{A})$ and p is sufficiently far from the boundary. This ensures that $\Delta := \min\{d(v, p) \mid v \in \text{conv}(\mathcal{A})\}$ is sufficiently greater than zero.

As can be seen in Table A.6 (Appendix A), all the algorithms have found a witness in no more than 4 iterations. Bear in mind that Theorem 2.6 states that the iteration complexity, in this case, does not depend on the tolerance ε , but on the constants Δ and R ; it is not hard to show that the ratio R/Δ approaches 1 as Δ increases. Figure 9 depicts the running times for this scenario.

Recall from Theorems 2.6 and 2.11 that the iteration complexity of TA and GT, when $p \notin \text{conv}(\mathcal{A})$, is

$$\mathcal{O}\left(\min\left\{\frac{R^2}{\Delta^2}, \frac{1}{c} \ln \frac{\delta_0}{\Delta}\right\}\right),$$

where c is from (10) and depends on the choice of pivot. Since GT looks for a pivot v that minimizes $v^T(p_k - p)$, usually $\sin \theta_k$ will be smaller for GT than TA (see Figure 3) and thus, we expect a larger constant c for GT than TA which explains the lower number of iterations for the former.

For the instances considered in this subsection, we observed $1.58 \leq R \leq 1.8$, $0.32 \leq \Delta \leq 0.41$ and $0.68 \leq \delta_0 \leq 0.81$. Moreover, by evaluating $\sin \theta_k$ throughout the iterations, we deduced that if the visibility factor c exists, its value must be less than 0.02.

In addition, ASFW and GT perform better than SPG. The gap in running time between ASFW (and GT) and SPG increases as the number of points increases. Although the number of iterations required by SPG is similar to the others, remember that its cost per iteration is higher due to the projection cost of $\mathcal{O}(n \log n)$.

We remark that in these instances p_0 , the point of \mathcal{A} closest to p , can be either v_q or v_ℓ . As argued in the previous section, if $p_0 = v_q$, then v_ℓ is a strict pivot and can be chosen by GT and ASFW in the first iteration. Since, for such instances, the projection of p onto $\text{conv}(\mathcal{A})$ is a convex combination of v_q and v_ℓ , these two algorithms can find it at the first iteration, explaining the results in Table A.6.

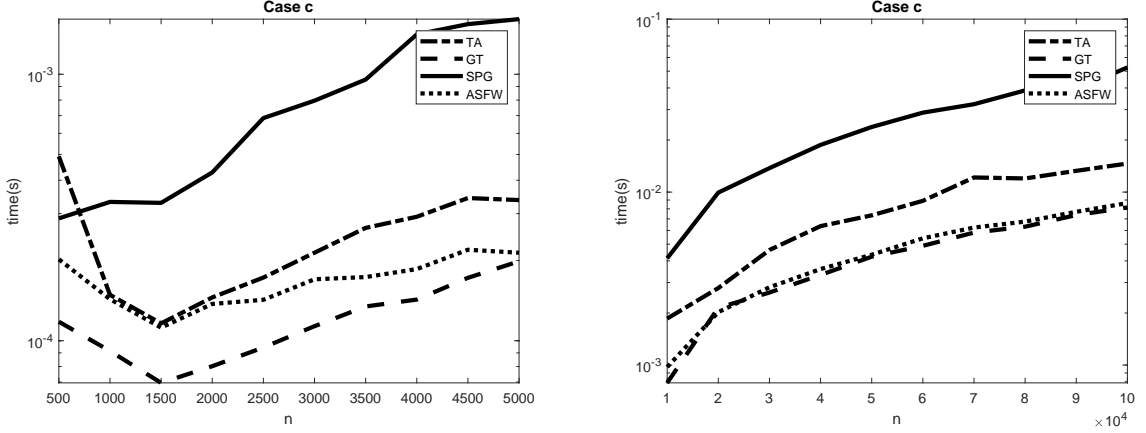


Figure 9: Case (c). Running times (in seconds) for dimension $m = 100$ and increasing n , the number of points in \mathcal{A} .

6.1.4. Case (d): $p \notin \text{conv}(\mathcal{A})$ with visibility factor close to zero.

Motivated again by Example 2.12, now we generate p as in case (b) but we multiply its coordinates by 1.01. Here we add v_s , as in case (b) and, as explained in Section 6.1.2, we expect to obtain instances with visibility factor close to zero. In fact, by analyzing the TA iterations for this set of instances, we observed that if a positive factor c exists its value will be less than $O(10^{-4})$. Moreover, for the instances generated in this case, we have $0.007 \leq \Delta \leq 0.008$, $1.34 \leq R \leq 1.54$, and $0.57 \leq \delta_0 \leq 0.75$. Thus, from Theorem 2.6 (and Theorem 2.11) in the worst case TA may spend $O(10^4)$ iterations to find a witness.

As can be seen in Table A.6, the average number of TA and GT iterations is very high in the comparison to the other algorithms and this translates to the running times shown in Figure 10. We remark that the number of SPG and ASFW iterations also increased in this case when compared to case (c). Note that the Triangle Algorithm takes practically the same time as GT.

Remarkably, the away directions enabled ASFW to avoid the zigzag phenomenon seen in the Example 2.12: it was the fastest algorithm, followed by SPG.

6.2. Linear Programming Feasibility Problem

The convex hull membership problem is also related to the linear programming feasibility problem. The feasible set of a linear programming problem can be given by $\Omega = \{x \in \mathbb{R}^n \mid Ax = b, x \geq 0\}$, with matrix $A = [a_1 \ a_2 \ \dots \ a_n]$, and vectors $b, a_i \in \mathbb{R}^m$. We are interested in elements of Ω with bounded norm, let us say $\|x\|_1 \leq N$. In other words, we would like to decide whether the intersection of Ω with the half-space defined by $H_N = \{x \in \mathbb{R}^n \mid e^T x \leq N\}$, for a given $N > 0$, is non-empty. This (more specific) LP feasibility problem can be cast as

$$\begin{pmatrix} A & 0 & -b \\ e^T & 1 & -N \\ 0^T & 0 & 1 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \\ \gamma \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \frac{1}{N+1} \end{pmatrix}, \quad (32)$$

$$e^T \alpha + \beta + \gamma = 1,$$

$$\alpha, \beta, \gamma \geq 0,$$

which is equivalent to the CHMP $p \in \text{conv}(\tilde{\mathcal{A}})$, where $p = (0^T, 0, \frac{1}{N+1})^T \in \mathbb{R}^{m+2}$ and the elements of $\tilde{\mathcal{A}}$ are the columns of the coefficient matrix in (32). It is not hard to show that $\Omega \cap H_N$ is nonempty if and only if $p \in \text{conv}(\tilde{\mathcal{A}})$.

Remark 6.2. Let $p_k \in \text{conv}(\tilde{\mathcal{A}})$ be given such that $d(p_k, p) = \|p_k - p\| \leq \varepsilon R$, where R is the diameter of $\text{conv}(\tilde{\mathcal{A}})$, from (32). Then, $\hat{x} := (\frac{\alpha_1}{\gamma}, \dots, \frac{\alpha_n}{\gamma})^T$ satisfies $\|A\hat{x} - b\| \leq \varepsilon R / \gamma =: \varepsilon'$, $|e^T \hat{x} + \beta / \gamma - N| \leq \varepsilon'$, $|\gamma - \frac{1}{N+1}| \leq \varepsilon R$ and $\hat{x} \geq 0$.

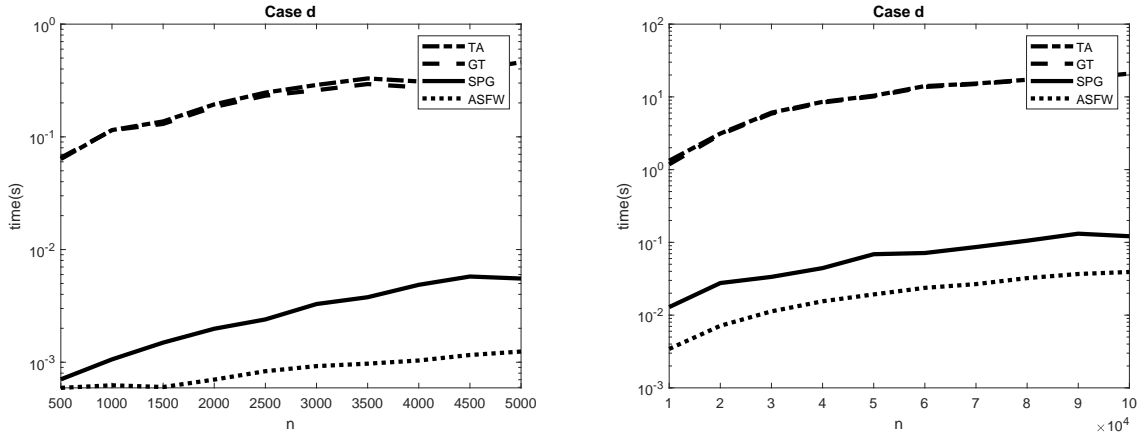


Figure 10: Case (d). Running times (in seconds) for dimension $m = 100$ and increasing n , the number of points in \mathcal{A} .

We consider instances of the LP feasibility problem

$$Ax = b, \quad x \geq 0 \quad \text{and} \quad e^T x \leq N, \quad (33)$$

and solve them as CHMPs by using TA, GT, SPG and ASFW, using the stopping criteria described in Section 5. We compared as well their performances against the Matlab LP solver `linprog` (which is an implementation of the Dual-Simplex algorithm) applied to (33).

For each instance, the columns of A were sampled from the uniform distribution on the unit sphere centered at $e \in \mathbb{R}^m$ (such that $a_{ij} \geq 0$). To generate feasible instances of (33), we select a random vector $x \in \mathbb{R}^n$, where each entry is sampled from the uniform distribution in $(0, 1)$, and compute b as $b = Ax$. Since the sum of n independent random variables with uniform distribution has expected value $n/2$ and variance $n/12$, and in these experiments we consider $n \leq 2000$, to ensure $e^T x \leq N$, we set $N = 1200$. For infeasible instances, we use the same procedure but at the end we multiply the first coordinate of b by -1 . In the following tables, for each pair m and n we considered 10 random instances.

For `linprog`, we test the feasibility of (33) considering a constraint violation tolerance of 10^{-3} (the largest value allowed by `linprog`) and we also disabled preprocessing procedures. For the other algorithms applied to the equivalent CHMP formulation, we use the same parameters of the experiments in Section 6.1 except for the tolerance ε , which here takes values in $\{10^{-6}, 10^{-7}\}$, and the maximum number of iterations, which is set to 10^6 . In the SPG implementation, based on preliminary experiments, we set $\lambda_{\min} = 10^{-10}$ and $\lambda_{\max} = 10^{10}$, and we use $M = 60$ for the nonmonotone line-search. It is worth mention that, for the instances we generated, the constant R , that depends on the problem geometry, is around 10^3 .

The running times and respective iterations count are shown in Tables 2, 3 and 4. In Table 2, we report time and iterations of TA, GT, ASFW and SPG for feasible problems, varying dimensions and tolerance ε . Table 3 shows the performance of `linprog` (for which the constraint violation tolerance is fixed).

Note that even GT and ASFW, which presented the best performances in the experiments of Section 6.1, struggled in these experiments, requiring many iterations and consequently presenting higher running times. This was somehow expected because Theorem 2.6 states the iteration complexity of TA and its variants as $O(\varepsilon^{-2})$ and now ε is three orders of magnitude smaller than in the previous section.

As for ASFW, we remark that the linear convergence rate appearing in equation (18) depends on μ , L , D , Ω_C and the dimension. In these experiments, $C := \text{conv } \tilde{\mathcal{A}}$ and the dimension is $m + 2$. The objective function of (4) yields that $\mu = L = 1$. From (32), we deduce that $D \geq N + 1$ and $\Omega_C \leq 1$. Therefore,

$$\frac{\Omega_C}{D(m+3)} \leq \frac{1}{(N+1)(m+3)}.$$

Table 2: CPU time (in seconds) and iterations count for feasible problems. Best times are in bold.

				TA		GT		ASFW		SPG	
		m	n	time	iter	time	iter	time	iter	time	iter
$\varepsilon = 10^{-6}$	50	200	0.5137	86491.4	0.3272	55952.9	0.1507	6286.7	0.0053	56.8	
	50	2000	0.2669	10318	0.0809	4472.4	0.0869	2190.3	0.2971	54.4	
	100	500	1.0028	74248.7	0.5599	46952.8	1.0131	27998.3	0.0083	122	
	200	2000	1.5370	29734.1	0.8193	17676.6	1.2437	18007.5	0.1833	268.5	
$\varepsilon = 10^{-7}$	50	200	2.1938	330132	1.2349	189164.6	0.7452	26773.2	0.0080	58.5	
	50	2000	0.6298	24588.7	0.1503	8173.1	0.2172	5847.9	0.5040	55.9	
	100	500	4.0366	313218.6	2.1494	186011.7	4.0959	112368.9	0.0148	125.2	
	200	2000	13.7279	245975.6	1.8809	39395.9	8.8099	123920.4	0.1923	266.2	

Table 3: Time (in seconds) for `linprog`

m	n	feasible	infeasible
50	200	0.0162	0.0103
50	2000	0.0294	0.0189
100	500	0.0297	0.0128
200	2000	0.2180	0.0424

Since, in our experiments, $N = 1200$ and $m = O(10^2)$, the above ratio is smaller than 10^{-5} and by replacing it in (18) we see that convergence factor is very close to 1, justifying the slow convergence of ASFW for these instances.

Table 4 shows the results (on average) for infeasible problems. We remark that in infeasible problems, apart from `linprog`, the algorithms stop when a witness is found and this depends only on the geometry of the problem, not on the tolerance ε . For the instances considered in this section we observed that R is around 10^3 and $0.14 \leq \Delta \leq 0.99$. We can see that the algorithm with the best performance with respect to time in feasible cases was SPG, followed by `linprog`. For infeasible problems, GT found a witness faster than the other benchmarked algorithms.

6.3. An example of CHMP in image classification

The MNIST database is a large database of handwritten digits, which is a very popular dataset for validating image classification algorithms (LeCun et al., 1998). It contains 60,000 samples of 28×28 grayscale images in its training set, and 10,000 images in its testing set. In the interesting study of Yousefzadeh (2021) it was shown that, for the MNIST dataset, all test points are considerably outside the convex hull of the set of training points. Furthermore, in such a case, the author suggested the following procedure for classification: given a testing point, we estimate its distance to the convex hull of each class³ and classify it according to the smallest of the distances. Yousefzadeh & Huang (2020) and Yousefzadeh (2021) reported an accuracy of 98.5% using this procedure.

³For MNIST, the classes are 0, 1, ..., 9.

Table 4: CPU time (in seconds) and iterations count for infeasible problems.

		TA		GT		ASFW		SPG	
m	n	time	iter	time	iter	time	iter	time	iter
50	200	0.0011	21.1	0.0003	13	0.0084	253.9	0.0018	3.5
50	2000	0.0016	35	0.0007	33.2	0.0015	33.2	0.0352	11.3
100	500	0.0014	56.9	0.0009	45.4	0.0015	45.4	0.0053	3.7
200	2000	0.0053	109.2	0.0045	106.8	0.0060	106.8	0.0513	7.6

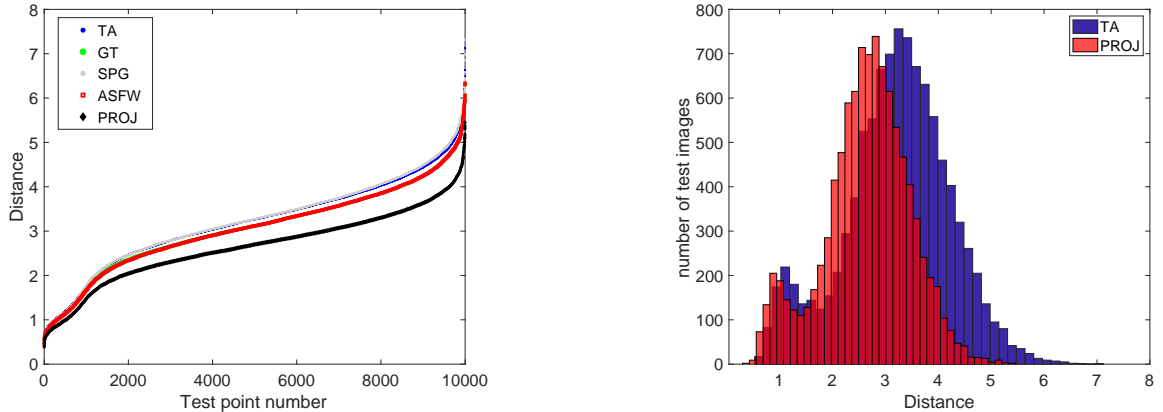


Figure 11: (Left) Distance estimates compared with the distance to the convex hull (PROJ). (Right) Histogram with the distances between the witness and the test point (TA), and the distances between the test point and its projection (PROJ).

Inspired by these studies, we apply the algorithms discussed in this work to obtain an estimate for the distance from a test point to the convex hull of (subclasses of) the training set. The difference from the work of [Yousefzadeh \(2021\)](#) is that we do not solve problem (3) exactly, but stop the algorithms as soon as a witness is found. Recall from (5) that the distance from a testing point to a witness is not greater than twice the distance from this test point to the convex hull. Thus, we wonder how the use of this estimate would impact the results/conclusions of [Yousefzadeh \(2021\)](#).

Firstly, we consider $\mathcal{A} \subset \mathbb{R}^{784}$ as the training set (with 60,000 vectorized images). Then, for each of the 10,000 test points, we solve the corresponding CHMPs with the same algorithms⁴ tested in the Section 6.1. As in ([Yousefzadeh, 2021](#)), we found that all elements of the test set are outside $\text{conv}(\mathcal{A})$. Similarly to Section 6.1.3, the fastest algorithms to obtain a witness for the 10,000 test points were GT and ASFW, which took approximately a total time of 12 and 13 minutes, respectively. The slowest was SPG which took 114 minutes. In order to obtain the exact distance from a testing point to the convex hull we also considered a classic version of SPG (without the stopping criteria of Section 5), which we call PROJ. PROJ took a total time of almost 9 hours to find the projection of all test points onto $\text{conv}(\mathcal{A})$.

To compare the distance estimates with the exact distances to the convex hull of the training set, we present in Figure 11 a graph with the distance (divided by 255 which is the largest pixel value) from all test points to their respective witnesses (depending on the method used) and a histogram to compare the (estimate) distance distributions for TA and PROJ. As expected, the distances between the witnesses and the test points do not exceed twice the distance to the projection.

Secondly, we apply the same procedure for classification as proposed in ([Yousefzadeh & Huang, 2020](#)), however we calculate the distance from the testing points to the witness in the convex hull of each class, and we classify the testing point according to the smallest of the distances. The best accuracy of 98% is achieved by TA in less than 10 minutes, which is not that far from the 98.5% accuracy of PROJ, that required more than 4 hours to process all test points.

7. Concluding remarks

We have shown that first-order methods, such as Frank-Wolfe-type and Projected Gradient-type, equipped with appropriate stopping criteria, are suitable for solving convex hull membership problems and are competitive with the recently proposed Triangle Algorithm, which is specific for CHMP.

By exploring the characterization of the so-called pivots, we showed that the Triangle Algorithm can be viewed as an inexact version of the Frank-Wolfe algorithm applied to CHMP whereas the latter, with stopping

⁴Here, $M = 3$ is used in the nonmonotone line search for SPG.

rules based on distance duality, lead us to a Greedy Triangle Algorithm. This algorithm coincides with an old one, due to von Neumann.

The proposed stopping criteria based on distance duality were essential for the first-order methods to handle the case $p \notin \text{conv}(\mathcal{A})$, avoiding the computation of the projection of the query point onto the convex hull, and also allowing for a fair numerical comparison with the Triangle Algorithm and its variants.

A comprehensive set of numerical experiments indicate which algorithm is preferable according to the geometry of the convex hull and the relative position of the query point. Overall, for the random instances considered in this paper, ASFW had a better performance, especially in “harder to solve” instances of CHMP, as those of Sections 6.1.2 and 6.1.4, where the classic versions of TA and FW usually present a zigzagging behavior.

We also gave examples of two potential applications of algorithms for CHMP, one in the linear programming feasibility problem and another in image classification problems. After casting LP feasibility problem as a CHMP, our experiments show that first-order methods can be superior to the classic dual-simplex algorithm and support the Greedy Triangle algorithm as an efficient candidate to detect infeasibility. The image classification experiment shows that the use of the distance from a query point to a witness instead of the distance from the point to the convex hull did not change considerably the classification accuracy, and it is about twenty times faster than computing the exact projection. A further computational study on other image classification datasets is subject of future investigation.

The performance of Frank-Wolfe-type algorithms in feasible instances of LP feasibility problems was different from what we observed in Sections 6.1.1 and 6.1.2. As discussed in Section 6.2, this has to do with the conditioning of the problem (which depends on the geometry of the convex hull), which makes things harder even for the linearly convergent ASFW. In future works we plan to investigate this precisely as well as possible acceleration strategies for FW-type methods applied to CHMP.

A deeper study on Spherical-TA (Kalantari & Zhang, 2022) and a comparison with the first-order methods considered in this paper shall be subject of our future investigations as well as extensions and applications of distance duality to semidefinite and conic programming (Kalantari, 2019b, 2020).

Acknowledgments

The authors are grateful to the anonymous referees for meaningful suggestions that helped to improve this paper. The authors want to acknowledge Brazilian agency CAPES (Coordenação para Aperfeiçoamento Pessoal do Ensino Superior) for financial support. **RF** thanks CAPES for the doctoral scholarship; **DG** thanks Brazilian agency CNPq (Conselho Nacional de Desenvolvimento Científico e Tecnológico) for grant 305213/2021-0; **LRS** thanks CNPq for grant 113190/2022-0.

References

- Awasthi, P., Kalantari, B., & Zhang, Y. (2020). Robust vertex enumeration for convex hulls in high dimensions. *Annals of Operations Research*, 295, 37–73. doi:10.1007/s10479-020-03557-0.
- Beck, A., & Shtern, S. (2017). Linearly convergent away-step conditional gradient for non-strongly convex functions. *Mathematical Programming*, 164, 1–27. doi:10.1007/s10107-016-1069-4.
- Birgin, E. G., Martínez, J. M., & Raydan, M. (2000). Nonmonotone Spectral Projected Gradient Methods on Convex Sets. *SIAM Journal on Optimization*, 10, 1196–1211. doi:10.1137/S1052623497330963.
- Birgin, E. G., Martínez, J. M., & Raydan, M. (2009). Spectral projected gradient methods. In C. A. Floudas, & P. M. Pardalos (Eds.), *Encyclopedia of Optimization* (pp. 3652–3659). Boston, MA: Springer US. doi:10.1007/978-0-387-74759-0_629.
- Birgin, E. G., Martínez, J. M., & Raydan, M. (2014). Spectral Projected Gradient Methods: Review and Perspectives. *Journal of Statistical Software*, 60. doi:10.18637/jss.v060.i03.
- Canon, M. D., & Cullum, C. D. (1968). A Tight Upper Bound on the Rate of Convergence of Frank-Wolfe Algorithm. *SIAM Journal on Control*, 6, 509–516. doi:10.1137/0306032.
- Clarkson, K. L. (2010). Coresets, sparse greedy approximation, and the Frank-Wolfe algorithm. *ACM Transactions on Algorithms*, 6, 1–30. doi:10.1145/1824777.1824783.
- Condat, L. (2016). Fast projection onto the simplex and the ℓ_1 ball. *Mathematical Programming*, 158, 575–585. doi:10.1007/s10107-015-0946-6.
- Dantzig, G. B. (1992). *An ϵ -Precise Feasible Solution to a Linear Program with a Convexity Constraint in $1/\epsilon^2$ Iterations Independent of Problem Size*. Technical Report SOL 92-5 Systems Optimization Laboratory, Dept of Management Science and Engineering, Stanford University Palo Alto, CA. URL: <https://stanford.edu/group/SOL/reports/SOL-92-5.pdf>.
- Epelman, M., & Freund, R. M. (2000). Condition number complexity of an elementary algorithm for computing a reliable solution of a conic linear system. *Mathematical Programming*, 88, 451–485. doi:10.1007/s101070000136.
- Frank, M., & Wolfe, P. (1956). An algorithm for quadratic programming. *Naval Research Logistics Quarterly*, 3, 95–110. doi:10.1002/nav.3800030109.

- Gonçalves, J. P., Storer, R. H., & Gondzio, J. (2009). A family of linear programming algorithms based on an algorithm by von Neumann. *Optimization Methods and Software*, 24, 461–478. doi:10.1080/10556780902797236.
- Grapiglia, G. N., & Sachs, E. W. (2017). On the worst-case evaluation complexity of non-monotone line search algorithms. *Computational Optimization and Applications*, 68, 555–577. doi:10.1007/s10589-017-9928-3.
- Grippo, L., Lampariello, F., & Lucidi, S. (1986). A Nonmonotone Line Search Technique for Newton’s Method. *SIAM Journal on Numerical Analysis*, 23, 707–716. doi:10.1137/0723046.
- Guélat, J., & Marcotte, P. (1986). Some comments on Wolfe’s ‘away step’. *Mathematical Programming*, 35, 110–119. doi:10.1007/BF01589445.
- Harman, R., & Lacko, V. (2010). On decompositional algorithms for uniform sampling from n -spheres and n -balls. *Journal of Multivariate Analysis*, 101, 2297–2304. doi:10.1016/j.jmva.2010.06.002.
- Jaggi, M. (2013). Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In S. Dasgupta, & D. McAllester (Eds.), *Proceedings of the 30th International Conference on Machine Learning* (pp. 427–435). Atlanta, Georgia, USA: PMLR volume 28 of *Proceedings of Machine Learning Research*. URL: <http://proceedings.mlr.press/v28/jaggi13.html>.
- Kalantari, B. (2014). A characterization theorem and an algorithm for a convex hull problem. *Annals of Operations Research*, 226, 301–349. doi:10.1007/s10479-014-1707-2.
- Kalantari, B. (2019a). An algorithmic separating hyperplane theorem and its applications. *Discrete Applied Mathematics*, 256, 59–82. doi:10.1016/j.dam.2018.05.009.
- Kalantari, B. (2019b). A Triangle Algorithm for Semidefinite Version of Convex Hull Membership Problem. *arXiv:1904.09854v2 [cs, math]*, . [arXiv:1904.09854v2](https://arxiv.org/abs/1904.09854v2).
- Kalantari, B. (2020). On the Equivalence of SDP Feasibility and a Convex Hull Relaxation for System of Quadratic Equations. *arXiv:1911.03989v2 [cs, math]*, . [arXiv:1911.03989v2](https://arxiv.org/abs/1911.03989v2).
- Kalantari, B., & Park, J. Y. (2014). Tree convex hull theorems on triangles and circles. *Honam Mathematical J.*, 36, 787–794. doi:10.5831/HMJ.2014.36.4.787.
- Kalantari, B., & Zhang, Y. (2022). Algorithm 1024: Spherical Triangle Algorithm: A Fast Oracle for Convex Hull Membership Queries. *ACM Transactions on Mathematical Software*, 48, 23:1–32. doi:10.1145/3516520.
- Lacoste-Julien, S., & Jaggi, M. (2015). On the global linear convergence of Frank-Wolfe optimization variants. In *Proceedings of the 28th International Conference on Neural Information Processing Systems* (pp. 496–504). Cambridge, MA, USA: MIT Press volume 1 of *NIPS’15*. doi:10.5555/2969239.2969295.
- LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86, 2278–2324. doi:10.1109/5.726791.
- Li, M., & Kalantari, B. (2013). Experimental study of the convex hull decision problem via a new geometric algorithm. In *23rd Annual Fall Workshop on Computational Geometry, City College of New York*.
- Peña, J., Rodríguez, D., & Soheili, N. (2016). On the von Neumann and Frank–Wolfe Algorithms with Away Steps. *SIAM Journal on Optimization*, 26, 499–512. doi:10.1137/15M1009937.
- Yousefzadeh, R. (2021). Deep Learning Generalization and the Convex Hull of Training Sets. *arXiv:2101.09849 [cs, math]*, . [arXiv:2101.09849](https://arxiv.org/abs/2101.09849).
- Yousefzadeh, R., & Huang, F. (2020). Using Wavelets and Spectral Methods to Study Patterns in Image-Classification Datasets. *arXiv:2006.09879 [cs, eess, math]*, . [arXiv:2006.09879](https://arxiv.org/abs/2006.09879).
- Zhang, H., & Hager, W. W. (2004). A Nonmonotone Line Search Technique and Its Application to Unconstrained Optimization. *SIAM Journal on Optimization*, 14, 1043–1056. doi:10.1137/S1052623403428208.
- Zhang, Y., & Kalantari, B. (2016). The Triangle Algorithm with Relaxed and Randomized Pivots. In *Proceedings of 26th Fall Workshop on Computational Geometry*. New York, NY: CUNY Graduate Center. https://matthewpjohnson.org/fwcg2016/FWCG_2016_paper_30.pdf.

Appendix A. Tables with numerical results from Section 6.1

Table A.5: Iterations count (average) for the case where $p \in \text{conv}(\mathcal{A})$

n	Case (a)				Case (b)	
	TA	ASFW	GT	SPG	ASFW	SPG
500	2557.3	573.9	662.2	23.7	12	8
1000	1544.8	242.8	247.9	15.9	12.5	8.8
1500	1428.1	195.0	196.5	15.0	13	9.3
2000	1373.0	167.4	169.7	13.8	12	8.0
2500	1317.2	153.4	153.7	13.1	12	7.8
3000	1345.9	147.3	146.1	13.1	12.5	10
3500	1288.7	139.1	139.5	13.1	12	8.6
4000	1295.5	132.3	133.3	13.0	12	8.9
4500	1270.0	129.8	129.8	13.0	13	10.3
5000	1309.7	126.1	126.6	13.0	12	8.4
1×10^4	1261.5	110.7	110.5	12.0	13	11.8
2×10^4	1254.6	99.0	99.0	11.9	13	12.4
3×10^4	1220.8	93.4	93.4	12.0	13	10.6
4×10^4	1251.6	88.0	87.8	12.0	12	11.0
5×10^4	1235.0	86.7	86.7	12.0	12	12.6
6×10^4	1243.1	84.1	84.1	12.0	13	12.8
7×10^4	1265.8	83.5	83.4	12.0	12	9.4
8×10^4	1247.3	82.2	82.2	12.0	13	13.1
9×10^4	1233.6	80.7	80.8	12.0	13	11.0
10^5	1225.5	79.3	79.4	12.0	12	10.2

Table A.6: Iterations count (average) for the cases where $p \notin \text{conv}(\mathcal{A})$

n	Case (c)				Case (d)			
	TA	ASFW	GT	SPG	TA	ASFW	GT	SPG
500	2.3	1	1	1.3	6570.6	9.2	6575.2	4.6
1000	3	1	1	1.3	7347.5	9.1	7358.4	4.0
1500	1.9	1	1	1.2	6474	9.2	6483.2	4.8
2000	2.6	1	1	1.3	7233.5	9.0	7246.2	4.5
2500	2.2	1	1	1.5	6604	9.1	6616.4	4.1
3000	2.2	1	1	1.5	6657.9	9.2	6670.2	4.6
3500	3	1	1	1.4	6519.3	9.2	6529.5	4.0
4000	2.0	1	1	1.7	5553.0	9.2	5560.4	4.8
4500	2.7	1	1	1.5	6275.2	9.2	6283.2	5.1
5000	2.3	1	1	1.5	6628.0	9.1	6638.5	4.3
1×10^4	3.3	1	1	1.3	6154.4	9.1	6165.8	4.9
2×10^4	2.1	1	1	1.8	5671.3	9.2	5680.0	5.2
3×10^4	2.6	1	1	1.6	6184.6	9.0	6201.9	4.0
4×10^4	2.9	1	1	1.6	6146.2	9.0	6159.5	4.0
5×10^4	2.6	1	1	1.6	5795.8	9.2	5811.4	5
6×10^4	2.6	1	1	1.6	6390.5	9.1	6416.2	4.4
7×10^4	3.4	1	1	1.5	5890.2	9.0	5903.4	4.6
8×10^4	2.7	1	1	1.6	5783.6	9.1	5803.0	4.9
9×10^4	2.6	1	1	1.4	5324	9.2	5334.4	5
10^5	2.6	1	1	1.7	5542.4	9.0	5560.3	4.4