# Avoiding bad steps in Frank-Wolfe variants

Francesco Rinaldi[*]        Damiano Zeffiro[†]

November 22, 2022

## Abstract

The study of Frank-Wolfe (FW) variants is often complicated by the presence of different kinds of "good" and "bad" steps. In this article, we aim to simplify the convergence analysis of specific variants by getting rid of such a distinction between steps, and to improve existing rates by ensuring a non-trivial bound at each iteration.

In order to do this, we define the Short Step Chain (SSC) procedure, which skips gradient computations in consecutive short steps until proper conditions are satisfied. This algorithmic tool allows us to give a unified analysis and converge rates in the general smooth non convex setting, as well as a linear convergence rate under a Kurdyka-Łojasiewicz (KL) property. While the KL setting has been widely studied for proximal gradient type methods, to our knowledge, it has never been analyzed before for the Frank-Wolfe variants considered in the paper.

An angle condition, ensuring that the directions selected by the methods have the steepest slope possible up to a constant, is used to carry out our analysis. We prove that such a condition is satisfied, when considering minimization problems over a polytope, by the away step Frank-Wolfe (AFW), the pairwise Frank-Wolfe (PFW), and the Frank-Wolfe method with in face directions (FDFW).

**Keywords:** Nonconvex optimization, First-order optimization, Frank-Wolfe variants, Kurdyka-Łojasiewicz property.
**MSC Classification:** 46N10, 65K05, 90C06, 90C25, 90C30

## 1   Introduction

The Frank-Wolfe method [25] and its variants (see, e.g., [26], [45] and references therein) provide a valid alternative to projected gradient approaches for the constrained optimization of a smooth objective $f : \mathbb{R}^n \to \mathbb{R}$, in settings where projecting on the feasible set may be unpractical. These methods have found many applications in sparse and structured optimization (see, e.g., [9], [26], [33], [37], [54] and references therein).

In this paper, we aim to overcome an annoying issue affecting the analysis of some FW variants, that is the presence of "bad iterations", i.e., iterations where we cannot

---

[*]Dipartimento di Matematica "Tullio Levi-Civita", Università di Padova, Italy (`rinaldi@math.unipd.it`)

[†]Dipartimento di Matematica "Tullio Levi-Civita", Università di Padova, Italy (`damiano.zeffiro@math.unipd.it`)

show good progress. This happens when we are forced to take a short step along the search direction to guarantee feasibility of the iterate. The number of short steps typically needs to be upper bounded in the convergence analysis with "ad hoc" arguments (see, e.g., [26] and [45]). The main idea behind our method is to chain several short steps by skipping gradient updates until proper conditions are met.

## 1.1 Related work

**FW variants.** The main drawback of the classic FW algorithm is its slow $O(1/k)$ convergence rate for convex objectives. This rate is tight even for strongly convex objectives on polytopes, due to a well understood zig-zagging behaviour near optima on the boundary (see, e.g., [22] and [65]). The study of assumptions and variants leading to faster rates is a rapidly developing field.

Alternative or modified directions moving away from "bad" vertices or atoms have a long history, starting at least with the work of Wolfe [65] (see [43] and [45] for recent references). In addition to considering new directions, the works [19] and [20] propose strategies to skip the linear minimization oracle (LMO) computation from time to time by caching linear minimizers, while the recent work [43] for optimization on polytopes applies recursively a FW variant to smaller polytopes. However, to our knowledge, no strategy to avoid short steps has been discussed in these previous works.

For smooth strongly convex objectives, the convergence rates of many of these "improved directions" FW variants is linear on polytopes (see, e.g., [8] and [45]). Furthermore, in [41] it was proved that convergence rate of an AFW variant is adaptive to Hölderian error bound conditions interpolating between the general convex case and the strongly convex one.

A different approach, adopted in the general smooth convex setting, is to use FW variants to approximate projections. In particular, the conditional gradient sliding method uses the FW method to approximate projections on the feasible set within a projected gradient scheme (see, e.g., [32] and [46]). Another approach introduced in [23] for smooth convex objectives implicitly uses the Non Negative Matching Pursuit (NNMP) algorithm to compute an approximate projection of the negative gradient on the tangent cone. To our knowledge, however, conditional gradient sliding approaches always lead to a sublinear $O(1/\varepsilon)$ LMO complexity, and the approach in [23] does not lead to any improvement on the $O(1/\varepsilon)$ worst case gradient complexity of the classic FW.

Outside the projection free setting, in [52] a procedure making multiple steps without updating the gradient (in a fashion similar to our SSC) is defined, and it is claimed that the approach traces the piecewise linear projection curve on polytopes, thus leading to the same linear convergence rate of the standard projected gradient method in the strongly convex setting.

In the non convex setting, for the classic FW algorithm a convergence rate of $O(1/\sqrt{k})$ was proved in [44] and then extended to other variants in [17] and [58].

**KL property.** The KL property (see, e.g., [4], [11] and [12]) has been extensively applied to compute the convergence rates of proximal subgradient type methods (see, e.g., [4], [5], [13], [64] and [66]). Furthermore, for convex objectives, it has been proved that Hölderian error bound conditions are a particular case of this property [13]. However, we are not aware of previous applications to the Frank-Wolfe variants under study in this paper.

**Angle condition.** The analysis of unconstrained descent methods often relies on some version of an angle condition, imposing an upper bound on the angle between the negative gradient and the descent direction selected by the method (see, e.g., [1], [29] and [67]). However, due to the presence of short steps and full FW steps, these analyses do not extend to our setting in a straightforward way.

In Section 3, we present an angle condition for optimization over a convex set. While to our knowledge this extension is novel for first order optimization methods, analogous conditions can be found in the context of direct search methods for linearly constrained derivative free optimization (see, e.g., [42] and [48]), imposed on the smallest angle between the negative gradient and a search direction. Finally, we remark that our condition was somehow used, but not stated explicitly, in [8] and [45] within the context of smooth strongly convex optimization over polytopes.

## 1.2 Contributions

Our main contributions are twofold:

- We formulate an angle condition for projection free methods, and prove that it leads to linear convergence in the number of "good steps" for non convex objectives satisfying a KL inequality. We show that this condition applies to the away step Frank-Wolfe (AFW), the pairwise Frank-Wolfe (PFW) and the FW method with in face directions (FDFW) (see, e.g., [26], [45], [28] and [31]) on polytopes. First, we give linear rates for good steps in Proposition 3.2. Then, we give global asymptotical rates under the assumption that the number of bad steps between two good steps is bounded in Proposition 3.3. We apply this result to FW variants in Corollary 3.1.

- We define the SSC procedure, which can be applied to all the FW variants listed in the first point, and show that it gets improvements on known rates (see Table 1 in Section 4). In particular, we prove that it leads to global linear convergence rates with no bad steps (see Lemma 4.3 and Corollary 4.3) under a global KL inequality and the angle condition. We then prove that we have local linear convergence rates and asymptotical linear convergence rates under a local KL property as well (see Theorem 4.2 and Corollary 4.2). This, to our knowledge, is the first (bad step free) linear convergence rate for FW variants under the KL inequality. In the general smooth non convex case, we further prove, under the angle condition, a $O(1/\sqrt{k})$ convergence rate with respect to a specific measure of non-stationarity for the iterates, that is the projection of the negative gradient on the convex cone of feasible directions (see Theorem 4.1, Corollary 4.1 and Remark 3).

While here we apply our framework only to the AFW, the PFW, and the FDFW on polytopes, we remark that our results hold for projection free methods on generic convex sets. In an extended version of this paper [60] we show applications on convex sets with smooth boundary for FW variants and methods using orthographic retractions (see also [2], [6], [47] and references therein).

The reasons why eliminating bad steps truly makes a difference in our context are the following:

- it rules out impractical convergence rates due to a large number of bad steps. An interesting example is given by the rate guarantee reported in [45] for the pairwise

Frank-Wolfe (PFW) variant on the $N-1$ dimensional simplex. This guarantee is indeed more loose than for the other variants, because there is no satisfactory bound on the number of such problematic steps (there is a best known bound of $3N!$ bad steps for each good step);

- it eliminates the dependence of the convergence rates on the support of the starting point (see, e.g., [35] and [43]). This dependence can significantly affect the performance of FW variants on smooth non convex optimization problems [24].

Finally, while beyond the scope of this paper, we mention that bad steps lead to a slow active set identification for the AFW, when compared to the "one shot" identification property characterizing proximal gradient methods and active set strategies (see [24], [53] and references therein). More precisely, analyses in recent works ([16], [17] and [27]) show that a number of bad steps equal to the number of "wrong" atoms is performed by the method in a sufficiently small neighborhood of a solution to identify its support.

## 1.3 Paper structure

The structure of the paper is as follows. In Section 2, we define some notation and state some preliminary results from convex analysis. In Section 3, we introduce the angle condition for first-order projection free methods, show examples of FW variants satisfying the condition and prove linear convergence in the number of good steps. We define the SSC procedure in Section 4, where we also state the main convergence results. Preliminary numerical results are reported in Section 6, while the missing proofs can be found in the appendix.

## 2 Notation and preliminaries

We consider the following constrained optimization problem:

$$\min \{f(x) \mid x \in \Omega\}. \tag{2.1}$$

In the rest of the article $\Omega$ is a compact and convex set and $f \in C^1(\Omega)$ with $L$-Lipschitz gradient:

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\| \text{ for all } x, y \in \Omega.$$

We define $D$ as the diameter of $\Omega$, $\hat{c} = c/\|c\|$ for $c \in \mathbb{R}^n/\{0\}$ and $\hat{c} = 0$ for $c = 0$. For sequences we write $\{x_k\}$ instead of $\{x_k\}_{k \in I}$ when $I$ is clear from the context, with $[j : i] = \{j, j+1, ..., i-1, i\}$. For $a, b \in \mathbb{R} \cup \{\pm\infty\}$ we denote as $[a < f(x) < b]$ the set $\{x \in \Omega \mid f(x) \in (a, b)\}$, with analogous definitions for non strict inequalities. For subsets $C, D$ of $\mathbb{R}^n$ we define $\text{dist}(C, D)$ as

$$\text{dist}(C, D) = \inf\{\|y - z\| \mid z \in C, \ y \in D\},$$

$B_R(C)$ as the neighborhood $\{x \in \mathbb{R}^n \mid \text{dist}(C, x) < R\}$ of $C$ of radius $R$ and in particular $B_R(x)$ as the open euclidean ball of radius $R$ and center $x$. When $C$ is closed and convex we define as $\pi(C, \cdot)$ the projection on $C$. If $C$ is a cone then we denote with $C^*$ its polar.

We now state some elementary properties related to the tangent and the normal cones, where for $\bar{x} \in \Omega$ we denote with $T_\Omega(\bar{x})$ and $N_\Omega(\bar{x})$ the tangent and the normal cone to $\Omega$ in $\bar{x}$ respectively. The next proposition (from [61], Theorem 6.9) characterizes these cones for closed convex subsets of $\mathbb{R}^n$.

**Proposition 2.1.** *Let $\Omega$ be a closed convex set. For every point $\bar{x} \in \Omega$ we have*

$$T_\Omega(\bar{x}) = \mathrm{cl}\{w \mid \exists \lambda > 0 \text{ with } \bar{x} + \lambda w \in \Omega\},$$
$$\mathrm{int}(T_\Omega(\bar{x})) = \{w \mid \exists \lambda > 0 \text{ with } \bar{x} + \lambda w \in \mathrm{int}(\Omega)\},$$
$$N_\Omega(\bar{x}) = T_\Omega(\bar{x})^* = \{v \in \mathbb{R}^n \mid (v, y - \bar{x}) \leq 0 \ \forall \ y \in \Omega\}.$$

We have the following formula connecting the supremum of a linear function "slope" along feasible directions to the tangent and the normal cone:

**Proposition 2.2.** *If $\Omega$ is a closed convex subset of $\mathbb{R}^n$, $\bar{x} \in \Omega$ then for every $g \in \mathbb{R}^n$*

$$\max\left\{0, \sup_{h \in \Omega \setminus \{\bar{x}\}} \left(g, \frac{h - \bar{x}}{\|h - \bar{x}\|}\right)\right\} = \mathrm{dist}(N_\Omega(\bar{x}), g) = \|\pi(T_\Omega(\bar{x}), g)\|.$$

This property is a consequence of the Moreau-Yosida decomposition [61] and we refer the reader to the Appendix for a detailed proof. On polytopes, a geometric interpretation is that the smallest angle between $g$ and a descent direction $d$ feasible in $\bar{x}$ is achieved for $d = \pi(T_\Omega(\bar{x}), g)$.

In the rest of the article to simplify notations we often use $\pi_{\bar{x}}(g)$ as a shorthand for $\|\pi(T_\Omega(\bar{x}), g)\|$. Then, by Proposition 2.2, first order stationarity conditions in $\bar{x}$ for the gradient $-g$ become equivalent to $\pi_{\bar{x}}(g) = 0$.

In the computation of the convergence rates, we often make the following assumption.

**Assumption 2.1.** *Given a stationary point $x^* \in \Omega$, there exists $\eta, \delta > 0$ such that for every $x \in [f(x^*) < f < f(x^*) + \eta] \cap B_\delta(x^*)$*

$$\pi_x(-\nabla f(x)) \geq \sqrt{2\mu}(f(x) - f(x^*))^{\frac{1}{2}}. \tag{2.2}$$

We refer the reader to the extended version [60] of this article for a study of convergence rates under a more general inequality, interpolating between (2.2) and the generic non convex case. Let now $i_\Omega$ be the indicator function of $\Omega$ so that $i_\Omega(x) = 0$ in $\Omega$ and $i_\Omega(x) = +\infty$ otherwise. It can easily be seen that (2.2) is a special case of the KL inequality (see, e.g., [4], [5] and [13]) with exponent $\frac{1}{2}$

$$\mathrm{dist}(0, \partial f_\Omega(x)) \geq \sqrt{2\mu}(f_\Omega(x) - f_\Omega(x^*))^{\frac{1}{2}} \tag{2.3}$$

for $f_\Omega = f + i_\Omega$, using that

$$\pi_x(-\nabla f(x)) = \mathrm{dist}(-\nabla f(x), N_\Omega(x)) = \mathrm{dist}(0, \partial(f + i_\Omega)(x)), \tag{2.4}$$

with the last equality following by Proposition 2.2. For convex objectives, condition (2.2) is therefore implied by the Holderian error bound $f(x) - f(x^*) \geq \gamma \mathrm{dist}(x, \mathcal{X}^*)^2$, for $\mathcal{X}^*$ set of solutions of Problem (2.1) (see [13, Corollary 6]), which in turn is implied by $\mu-$ strong convexity (see, e.g., [40]). Under suitable assumptions (see Proposition 8.1) our KL condition is also implied by the classic Polyak-Lojasiewicz inequality $\|\nabla f(x)\| \geq \sqrt{2\mu}(f(x) - f(x^*))^{\frac{1}{2}}$ (from [50] and [57]). Finally, Assumption 2.1 is implied by the Luo Tseng error bound [51] under some mild separability conditions for stationary points (see [49, Theorem 4.1]). This error bound is known to hold in a variety of convex and non convex settings (see Section 5 and references in [49]).

# 3 An angle condition

Let $\mathcal{A}$ be a first-order optimization method defined for smooth functions on a closed subset $\Omega$ of $\mathbb{R}^n$. We assume that given first-order information $(x_k, \nabla f(x_k))$ the method always selects $x_{k+1}$ along a feasible descent direction, so that for $(x, g) \in \Omega \times \mathbb{R}^n$ we can define

$$\mathcal{A}(x, g) \subset T_\Omega(x) \cap \{y \in \mathbb{R}^n \mid \langle g, y \rangle > 0\} \cup \{0\}$$

as the possible descent directions selected by $\mathcal{A}$ when $x = x_k$, $g = -\nabla f(x_k)$ for some $k$ (see Algorithm 1). When $x$ is first-order stationary, we set $\mathcal{A}(x, g) = \{0\}$, otherwise we always assume $0 \notin \mathcal{A}(x, g) \neq \emptyset$.

---

**Algorithm 1:** First-order method

---

**Initialization.** $x_0 \in \Omega$, $k := 0$.
1. If $x_k$ is stationary, then STOP
2. select a descent direction $d_k \in \mathcal{A}(x_k, -\nabla f(x_k))$
3. set $x_{k+1} = x_k + \alpha_k d_k$ for some stepsize $\alpha_k \in [0, \alpha_k^{\max}]$
4. set $k := k + 1$, go to Step 1.

---

We want to formulate an angle condition for the descent directions selected by $\mathcal{A}$, with respect to the infimum of the angles achieved with feasible descent directions. In order to do that, we define the directional slope lower bound as

$$\mathrm{DSB}_{\mathcal{A}}(\Omega, x, g) = \inf_{d \in \mathcal{A}(x, g)} \frac{\langle g, d \rangle}{\pi_x(g) \|d\|}$$

if $0 \notin \mathcal{A}(x, g)$. Otherwise $x$ is stationary for $-g$, $\pi_x(g) = 0$ and we set $\mathrm{DSB}_{\mathcal{A}}(\Omega, x, g) = 1$. Then with this definition it immediately follows $\mathrm{DSB}_{\mathcal{A}}(\Omega, x, g) \leq 1$ by Proposition 2.2. Notice also that when $x \in \mathrm{int}(\Omega)$ then $\mathrm{DSB}_{\mathcal{A}}(\Omega, x, g)$ is simply a lower bound on $\cos(\theta_{g,d})$ with $\theta$ the angle between $g$ and a descent direction $d$:

$$\mathrm{DSB}_{\mathcal{A}}(\Omega, x, g) = \inf_{d \in \mathcal{A}(x, g)} \frac{\langle g, d \rangle}{\|g\| \|d\|} \tag{3.1}$$

and thus imposing $\mathrm{DSB}_{\mathcal{A}}(\Omega, x, g) \geq \tau$ we retrieve the angle condition [1, equation (20)]. We remark that the RHS of (3.1) defining the unconstrained angle condition is also considered in the constrained setting in [23] (referred to as alignment condition), as a tool to evaluate potential descent directions. However, without $\pi_x(g)$ in the denominator no uniform lower bound can be given for the RHS, and therefore no worst case linear convergence rate (the rate given in [23, Corollary 3.6] is in fact $O(1/k)$). Given a subset $P$ of $\Omega$ we can finally define the slope lower bound

$$\mathrm{SB}_{\mathcal{A}}(\Omega, P) = \inf_{\substack{g \in \mathbb{R}^n \\ x \in P}} \mathrm{DSB}_{\mathcal{A}}(\Omega, x, g) = \inf_{\substack{g : \pi_x(g) \neq 0 \\ x \in P}} \mathrm{DSB}_{\mathcal{A}}(\Omega, x, g).$$

For simplicity if $P = \Omega$ we write $\mathrm{SB}_{\mathcal{A}}(\Omega)$ instead of $\mathrm{SB}_{\mathcal{A}}(\Omega, \Omega)$.

We now show a few examples of Frank-Wolfe variants satisfying the following *angle condition*

$$\text{SB}_\mathcal{A}(\Omega) = \tau > 0, \tag{3.2}$$

i.e. cases where the slope lower bound is strictly greater than 0.

## 3.1 Frank-Wolfe variants over polytopes and the angle condition

We now consider the AFW, PFW and FDFW and show that the angle condition is satisfied when $\Omega$ is a polytope. The AFW and PFW depend on a set of "elementary atoms" $A$ such that $\Omega = \text{conv}(A)$. Given $A$, for a base point $x \in \Omega$ we can define

$$S_x = \{S \subset A \mid x \text{ is a proper convex combination of all the elements in } S\},$$

the family of possible active sets for $x$. In the rest of the article $A$ is always clear from the context and for simplicity we write PFW, AFW instead of $\text{PFW}_A$, $\text{AFW}_A$. For $x \in \Omega$, $S \in S_x$, $d^{\text{PFW}}$ is a PFW direction with respect to the active set $S$ and gradient $-g$ iff

$$d^{\text{PFW}} = s - q \text{ with } s \in \text{argmax}_{s \in \Omega} \langle s, g \rangle \text{ and } q \in \text{argmin}_{q \in S} \langle q, g \rangle. \tag{3.3}$$

Similarly, given $x \in \Omega$, $S \in S_x$, $d^{\text{AFW}}$ is an AFW direction with respect to the active set $S$ and gradient $-g$ iff

$$d^{\text{AFW}} \in \text{argmax}\{\langle g, d \rangle \mid d \in \{d^{\text{FW}}, d^{AS}\}\}, \tag{3.4}$$

where $d^{\text{FW}}$ is a classic Frank-Wolfe direction

$$d^{\text{FW}} = s - x \text{ with } s \in \text{argmax}_{s \in \Omega} \langle s, g \rangle, \tag{3.5}$$

and $d^{\text{AS}}$ is the away direction

$$d^{\text{AS}} = x - q \text{ with } q \in \text{argmin}_{q \in S} \langle q, g \rangle. \tag{3.6}$$

The FDFW from [26], [31] (sometimes referred to as Decomposition invariant Conditional Gradient (DiCG) when applied to polytopes [28], [7]) relies only on the current point $x$ and the current gradient $-g$ to choose a descent direction and, unlike the AFW and the PFW, does not need to keep track of the active set.

The in face direction is defined as

$$d^F = x_k - x_F \text{ with } x_F \in \text{argmin}\{\langle g, y \rangle \mid y \in \mathcal{F}(x)\}$$

for $\mathcal{F}(x)$ the minimal face of $\Omega$ containing $x$. The selection criterion is then analogous to the one used by the AFW:

$$d^{\text{FD}} \in \text{argmax}\{\langle g, d \rangle \mid d \in \{d^F, d^{\text{FW}}\}\}. \tag{3.7}$$

We write $\text{SB}_{\text{FD}}, \text{DSB}_{\text{FD}}$ instead of $\text{SB}_{\text{FDFW}}, \text{DSB}_{\text{FDFW}}$ in the rest of the paper. When $\Omega$ is a polytope and $|A| < \infty$, the angle condition holds for the directions and the related FW variants we introduced. Before stating a lower bound for $\text{SB}_\mathcal{A}(\Omega)$ in this setting we need to recall the pyramidal width constant $\text{PWidth}(A)$ introduced in [45]. We refer

the reader to [59] and references therein for a discussion of various properties of this and related parameters.

We use here a characterization of $\text{PWidth}(A)$ proved in [55]:

$$\text{PWidth}(A) = \min_{\mathcal{F} \in \text{pfaces}(\Omega)} \text{dist}(\mathcal{F}, \text{conv}(A \setminus \mathcal{F})), \tag{3.8}$$

with $\text{pfaces}(\Omega)$ the set of proper faces of $\Omega$. We now introduce one key property of $\text{PWidth}(A)$ which relates it to the angle along the PFW direction. While we give a self contained proof of the lemma relying only on (3.8), we remark that the lemma can also be proved using [45, Theorem 3].

**Lemma 3.1.** *We have the following lower bound*

$$\frac{\langle g, d^{\text{PFW}} \rangle}{\|\pi(T_\Omega(x), g)\|} \geq \text{PWidth}(A).$$

*Proof.* We use $s, q$ and $S$ as in (3.3). For $z$ in $\Omega$ and $d$ feasible direction in $z$ we define as $\hat{\alpha}^{\max}(z, d)$ the maximal feasible stepsize in the direction $d$. Let $p = \pi(T_\Omega(x), g)$, and let $y$ be a maximizer of $\hat{\alpha}^{\max}(y, p)$ for $y \in S$. We have

$$\langle g, d^{\text{PFW}} \rangle = \langle g, (s - y) + (y - q) \rangle \geq \langle g, s - y \rangle \geq \langle g, (y + \hat{\alpha}^{\max}(y, p)p) - y \rangle$$
$$\geq \frac{\text{PWidth}(A)}{\|p\|} \langle g, p \rangle = \text{PWidth}(A)\|p\|, \tag{3.9}$$

where we used Lemma 8.2 in the third inequality, and $\langle g, p \rangle = \|p\|^2$ as it follows by the Moreau-Yosida decomposition in the last equality. $\square$

In order to define an angle condition for the FDFW, we use the following upper bound on $\text{PWidth}(A)$, independent from the particular set $A$ chosen to represent $\Omega$:

$$\text{PFWidth}(\Omega) = \min_{\substack{\mathcal{F}_1, \mathcal{F}_2 \in \text{pfaces}(\Omega) \\ \mathcal{F}_1 \cap \mathcal{F}_2 = \emptyset}} \text{dist}(\mathcal{F}_1, \mathcal{F}_2). \tag{3.10}$$

**Proposition 3.1.** $\text{SB}_{\text{PFW}}(\Omega) \geq \tau_p := \frac{\text{PWidth}(A)}{D}, \text{SB}_{\text{AFW}}(\Omega) \geq \frac{\tau_p}{2}, \text{SB}_{\text{FD}}(\Omega) \geq \frac{\tau_v}{2} := \frac{\text{PFWidth}(\Omega)}{2D}.$

*Proof.* Let $g$ be such that $\pi_x(g) \neq 0$. We have

$$\text{DSB}_{\text{PFW}}(\Omega, x, g) = \inf_{d^{\text{PFW}} \in \text{PFW}(x, g)} \frac{\langle g, d^{\text{PFW}} \rangle}{\|d^{\text{PFW}}\| \|\pi(T_\Omega(x), g)\|}$$
$$\geq \frac{\langle g, d^{\text{PFW}} \rangle}{D\|\pi(T_\Omega(x), g)\|} \geq \frac{\text{PWidth}(A)}{D},$$

where we used Lemma 3.1 in the last inequality.

Hence $\text{SB}_{\text{PFW}}(\Omega) \geq \frac{\text{PWidth}(A)}{D}$ follows by taking the inf on the LHS for $x \in \Omega$ and $g$ such that $\pi_x(g) \neq 0$ in (3.1). The inequality $\text{SB}_{\text{AFW}}(\Omega) \geq \frac{\text{PWidth}(A)}{2D}$ is a corollary since

$$\langle g, d^{\text{AFW}} \rangle \geq \frac{1}{2} \langle g, d^{\text{PFW}} \rangle,$$

8

as it follows immediately from the definitions (see also [45, equation (6)]).

The angle condition for the FDFW can be proved analogously to the angle condition for the AFW, where in Lemma 8.2 the RHS can be improved with $\text{PFWidth}(\Omega)$ instead of $\text{PWidth}(A)$ using that the active set $A'$ can be taken as the set of vertices of a face. $\quad\square$

**Remark 1.** *Results analogous to the ones in Proposition 3.1 can be proven relatively to the vertex facial distance* $\text{vf}(\Omega)$ *from [8]. More precisely, assuming* $A = V(\Omega)$, *for* $V(\Omega)$ *set of vertices of* $\Omega$, *and that the AFW and the PFW keep active sets of size at most* $\bar{s}$, *we have* $\text{SB}_{\text{PFW}}(\Omega) \geq \frac{\text{vf}(\Omega)}{\bar{s}D}$, $\text{SB}_{\text{AFW}}(\Omega) \geq \frac{\text{vf}(\Omega)}{2\bar{s}D}$ *as a consequence of [8, Lemma 3.1]. Furthermore, for the FDFW we have* $\text{SB}_{\text{FD}}(\Omega,\Omega_{\bar{s}}) \geq \frac{\text{vf}(\Omega)}{2\bar{s}D}$, *with* $x \in \Omega_{\bar{s}} \subset \Omega$ *iff there exists* $S \in S_x$ *such that* $|S| \leq \bar{s}$.

## 3.2 Linear convergence for good steps under the angle condition

Consider now a method following the scheme described by Algorithm 1 and with stepsize given by

$$\alpha_k = \min\left(\bar{\alpha}_k, \alpha_k^{\max}\right), \tag{3.11}$$

where

$$\bar{\alpha}_k = \frac{\langle -\nabla f(x_k), d_k \rangle}{L\|d_k\|^2}. \tag{3.12}$$

We notice that $\bar{\alpha}_k$ in (3.12) is a standard stepsize, often used in numerical tests with a properly tuned estimate for $L$ (see, e.g., [56]). The following lemma shows that at every iteration a sufficient decrease condition is satisfied, independently from the method $\mathcal{A}$, when using stepsize (3.12).

**Lemma 3.2.** *If* $\alpha_k \leq \bar{\alpha}_k$, *thus in particular for the stepsize* (3.11), *we have:*

$$f(x_k) - f(x_{k+1}) \geq \frac{L}{2}\|x_k - x_{k+1}\|^2. \tag{3.13}$$

The proof is straightforward and we defer it to the appendix.

Assume now that the method $\mathcal{A}$ used by Algorithm 1 satisfies the angle condition (3.2). We say that the algorithm performs a *full FW step* if

$$x_{k+1} \in \text{argmin}_{x \in \Omega} \langle \nabla f(x_k), x \rangle. \tag{3.14}$$

In the following proposition, we prove a general linear convergence rate in the number of *good steps*, i.e., the steps satisfying $\alpha_k = \bar{\alpha}_k$ or (3.14), under the assumption that the method $\mathcal{A}$ satisfies the angle condition (3.2), and that the KL inequality (2.2) holds for the objective function $f$ in Problem (2.1).

**Proposition 3.2.** *Let us assume that* $\mathcal{A}$ *satisfies the angle condition* (3.2), *and the objective function* $f$ *in Problem* (2.1) *satisfies condition* (2.2) *in* $x_k$ *and* $x_{k+1}$.

- *If* $\alpha_k = \bar{\alpha}_k$ *then*

$$f(x_{k+1}) - f(x^*) \leq \left(1 - \frac{\mu}{L}\tau^2\right)(f(x_k) - f(x^*)). \tag{3.15}$$

- *If the step* $k$ *is a full FW step then*

$$f(x_{k+1}) - f(x^*) \leq \left(1 + \frac{\mu}{L}\right)^{-1}(f(x_k) - f(x^*)). \tag{3.16}$$

9

*Proof.* Let $p_k = \|\pi(T_\Omega(x_{k+1}), -\nabla f(x_{k+1}))\|$ and $\tilde{p}_k = \|\pi(T_\Omega(x_{k+1}), -\nabla f(x_k))\|$. We have

$$
\begin{aligned}
|p_k - \tilde{p}_k| &= |\|\pi(T_\Omega(x_{k+1}), -\nabla f(x_{k+1}))\| - \|\pi(T_\Omega(x_{k+1}), -\nabla f(x_k))\|| \\
&\leq \|-\nabla f(x_{k+1}) + \nabla f(x_k)\| \leq L\|x_{k+1} - x_k\|,
\end{aligned}
\tag{3.17}
$$

where we used the 1-Lipschitzianity of projections in the first inequality.

If $\alpha_k = \bar{\alpha}_k$ then

$$
\begin{aligned}
f(x_{k+1}) &= f(x_k + \bar{\alpha}_k d_k) \leq f(x_k) - \frac{1}{2L}\left(\frac{\langle \nabla f(x_k), d_k \rangle}{\|d_k\|}\right)^2 \leq f(x_k) - \frac{\tau^2}{2L}p_{k-1}^2 \\
&\leq f(x_k) - \frac{\mu \tau^2}{L}(f(x_k) - f(x^*)),
\end{aligned}
\tag{3.18}
$$

where we used (8.3) in the first inequality, $\mathrm{SB}_{\mathcal{A}}^f(\Omega) = \tau$ in the second one, and condition (2.2) in the third one.

If the step $k$ is a full FW step then $\tilde{p}_k = 0$ because $x_{k+1} \in \mathrm{argmin}_{y \in \Omega}\langle \nabla f(x_k), y \rangle \Leftrightarrow -\nabla f(x_k) \in N_\Omega(x_{k+1}) \Leftrightarrow \|\pi(T_\Omega(x_{k+1}), -\nabla f(x_k))\| = 0$, where the last equivalence is true by Proposition 2.2. Then

$$
f(x_{k+1}) - f(x^*) \leq \frac{p_k^2}{2\mu} \leq \frac{(\tilde{p}_k + L\|x_{k+1} - x_k\|)^2}{2\mu} = \frac{L^2}{2\mu}\|x_{k+1} - x_k\|^2 \leq \frac{L}{\mu}(f(x_k) - f(x_{k+1})),
\tag{3.19}
$$

where we used (2.2) in the first inequality, (3.17) in the second, $\tilde{p}_k = 0$ and (8.4) in the last inequality. Then (8.4) and (3.16) follow by rearranging (3.18) and (3.19) respectively. $\qquad\square$

We finally report an asymptotic rate under the additional assumption that bad steps between two good steps are limited.

**Proposition 3.3.** *Assume that the number of bad steps between two good steps is limited and that $\mathcal{A}$ satisfies the angle condition (3.2). Then:*

- *every accumulation point of $\{x_k\}$ is stationary, and $f(x_k)$ is decreasing and convergent to $f^* \in \mathbb{R}$;*

- *if Assumption 2.1 holds for every stationary point in the level set $[f(x) = f^*]$, we have the asymptotic convergence rate:*

$$
f(x_k) - f(x^*) \leq M\bar{q}^{\bar{\gamma}(k)},
\tag{3.20}
$$

*for some $M > 0$, $\bar{\gamma}(k)$ number of good steps among the first $k$ steps and*

$$
\bar{q} = \max\left(\left(1 + \frac{\mu}{L}\right)^{-1}, \left(1 - \frac{\mu}{L}\tau^2\right)\right).
\tag{3.21}
$$

*Proof.* Let $k(j)$ be the subsequence of iterates associated to good steps, so that by assumption $k(j+1) - k(j)$ is bounded, and define $\tilde{k}(j) = k(j) - 1$ if $\alpha_{k(j)} = \bar{\alpha}_{k(j)}$, $\tilde{k}(j) = k(j)$ otherwise. Notice that $\tilde{k}(j+1) - \tilde{k}(j)$ is also bounded. By (8.4) we have that $\{f(x_k)\}$ is decreasing and thus convergent to $f^* \in \mathbb{R}$, and also that $\|x_k - x_{k+1}\| \to 0$.

10

With the notation used in Proposition 3.2 we now claim $p_{\tilde{k}(j)} \to 0$. In fact if $\alpha_{k(j)} = \bar{\alpha}_{k(j)}$ then

$$p_{\tilde{k}(j)}^2 = p_{k-1}^2 \le \frac{2L}{\tau^2}(f(x_k) - f(x_{k+1})) \to 0, \tag{3.22}$$

where we used (3.18) in the inequality, and if $k(j)$ is a full FW step then

$$p_{\tilde{k}(j)} \le p_{k(j)} \le \tilde{p}_{k(j)} + L\|x_{k(j)+1} - x_{k(j)}\| = L\|x_{k(j)+1} - x_{k(j)}\| \to 0, \tag{3.23}$$

where we used (3.17) in the first inequality and $\tilde{p}_{k(j)} = 0$ in the equality.

We therefore have $p_{\tilde{k}(j)} \to 0$. Equivalently, thanks to (2.4) we have $\mathrm{dist}(0, \partial f_\Omega(x_{\tilde{k}(j)})) \to 0$, so if $x^*$ is a limit point of $x_{\tilde{k}(j)}$ by lower semicontinuity of the subdifferential we must have $0 \in \partial f_\Omega(x^*)$, i.e., $x^*$ is stationary. In particular, by compactness $\{x_{\tilde{k}(j)}\}$ must converge to the set of stationary points. By the boundedness of $\|x_{k+1} - x_k\|$ and $\tilde{k}(j+1) - \tilde{k}(j)$ we also have that the set of limit points of $\{x_k\}$ coincides with the set of limit points of $\{x_{\tilde{k}(j)}\}$, and in particular it is a subset of stationary points contained in $[f(x) = f^*]$.

Let $\bar{\Omega} \subset [f(x) = f^*]$ be the set of limit points of $\{x_k\}$. By compactness (see [14, Lemma 6]), we have that for some fixed $\varepsilon, \eta > 0$, the KL property holds for every $x^* \in \bar{\Omega}$ with parameters $\varepsilon$ and $\eta$. Then for $k$ large enough $x_k \in B_\delta(x^*) \cap [f(x^*) < f < f(x^*) + \eta]$ for some $x^* \in \Omega$, and the asymptotic rates follow by Proposition 3.2. $\square$

For the three FW variants described before we can now give an asymptotic linear convergence rate in the number of good steps. We refer the reader to Table 1 for bounds on this number.

**Corollary 3.1.** *Let us assume that the objective function $f$ satisfies Assumption 2.1 for every stationary point in the level set $[f(x) = f^*]$ and $\Omega = \mathrm{conv}(A)$ with $|A| < +\infty$ in Problem* (2.1). *Then the AFW, the PFW and the FDFW converge at a rate*

$$f(x_k) - f(x^*) \le M\bar{q}_{gs}^{\bar{\gamma}(k)}, \tag{3.24}$$

*for some $M > 0$, with $\bar{\gamma}(k)$ the number of good steps among the first $k$ steps,*

$$\bar{q}_{gs} = \max\left(1 - \frac{\mu}{L}\left(\frac{\mathrm{PWidth}(A)}{2D}\right)^2, \left(1 + \frac{\mu}{L}\right)^{-1}\right) \tag{3.25}$$

*for the AFW,*

$$\bar{q}_{gs} = 1 - \frac{\mu}{L}\left(\frac{\mathrm{PWidth}(A)}{D}\right)^2 \tag{3.26}$$

*for the PFW, and*

$$\bar{q}_{gs} = \max\left(1 - \frac{\mu}{L}\left(\frac{\mathrm{PFWidth}(\Omega)}{2D}\right)^2, \left(1 + \frac{\mu}{L}\right)^{-1}\right) \tag{3.27}$$

*for the FDFW.*

*Proof.* For the AFW and the FDFW the rates (3.25) and (3.27) for good steps follow directly from (3.15) and (3.16) together with the bound on $\tau$ given in Proposition 3.1. Since the PFW never performs full FW steps, its rate (3.26) for good steps follow directly from (3.15) together with the bound on $\tau$ given in Proposition 3.1. Finally, given that the number of bad steps between two good steps is limited for all these methods (see [45, 43]), we have all the assumptions to apply Proposition 3.3. $\square$

11

# 4 First order projection free methods with SSC procedure

We introduce here the SSC procedure, and prove convergence rates both under the KL inequality (2.2) and in the generic non convex case.

## 4.1 The SSC procedure

The SSC procedure chains consecutive short steps, thus skipping updates for the gradient (and possibly for related information, like linear minimizers), until proper stopping conditions are met. Such a procedure, whose detailed scheme is given in Algorithm 3, can be easily embedded in a first-order approach (see Algorithm 2).

---

**Algorithm 2:** First-order method with SSC

---

**Initialization.** $x_0 \in \Omega$, $k = 0$.
1. **while** $x_k$ is not stationary:
2. $\quad g = -\nabla f(x_k)$
3. $\quad x_{k+1} = \mathrm{SSC}(x_k, g)$
5. $\quad k = k + 1$.

---

**Algorithm 3:** $\mathrm{SSC}(\bar{x}, g)$

---

**Initialization.** $y_0 = \bar{x}$, $j = 0$.
**Phase I**
1. $\quad$ select $d_j \in \mathcal{A}(y_j, g)$, $\alpha_{\max}^{(j)} \in \alpha_{\max}(y_j, d_j)$
2. $\quad$ **if** $d_j = 0$ **then:**
3. $\quad\quad$ **return** $y_j$
**Phase II**
4. $\quad$ compute $\beta_j$ with (4.2)
5. $\quad$ let $\alpha_j = \min(\alpha_{\max}^{(j)}, \beta_j)$
6. $\quad$ $y_{j+1} = y_j + \alpha_j d_j$
7. $\quad$ **if** $\alpha_j = \beta_j$ **then**:
8. $\quad\quad$ **return** $y_{j+1}$
9. $\quad$ $j = j + 1$, go to Step 1.

---

Given that the gradient $-g$ is constant during the SSC, this procedure is an application of $\mathcal{A}$ for the minimization of the linearized objective $f_g(z) = \langle -g, z - \bar{x} \rangle + f(\bar{x})$ with peculiar stepsizes and stopping criterion. More specifically, after a stationarity check (Phase I), the stepsize $\alpha_j$ is computed by taking the minimum between the maximal stepsize $\alpha_{\max}^{(j)}$ (which we always assume to be greater than 0) and an auxiliary stepsize $\beta_j$. The point $y_{j+1}$ generated in Phase II is always feasible since $\alpha_j \le \alpha_{\max}^{(j)}$ is always

smaller than the maximal feasible stepsize along the direction $d_j$. Notice that if the method $\mathcal{A}$ used in the SSC performs a FW step (see equation (3.5) for the definition of FW step), then the SSC terminates, with $\alpha_j = \beta_j$ or with $y_{j+1}$ global minimizer of $f_g$.

The auxiliary step size $\beta_j$ is defined as the maximal feasible stepsize for the trust region

$$\Omega_j = \bar{B}_{\|g\|/2L}(\bar{x} + \frac{g}{2L}) \cap \bar{B}_{\langle g, \hat{d}_j \rangle / L}(\bar{x}) \tag{4.1}$$

when $y_j \in \Omega_j$, otherwise the method stops returning $y_j$. Summarizing,

$$\beta_j = \begin{cases} 0 & \text{if } y_j \notin \Omega_j\,, \\ \beta_{\max}(\Omega_j, y_j, d_j) & \text{if } y_j \in \Omega_j\,, \end{cases} \tag{4.2}$$

where $\beta_{\max}(\Omega_j, y_j, d_j) = \max\{\beta \in \mathbb{R}_{\geq 0} \mid y_j + \beta d_j \in \Omega_j\}$ is the maximal feasible stepsize in the direction $d_j$ starting from $y_j$ with respect to $\Omega_j$. Since $\Omega_j$ is the intersection of two balls there is a simple closed form expression for $\beta_j$. In particular, using that $y_0 = \bar{x}$, if $d_0 \neq 0$ we have

$$\beta_0 = \frac{\langle g, \hat{d}_0 \rangle}{L \|d_0\|}\,,$$

which corresponds to (3.11) in the non maximal case, and where $\beta_0 > 0$ since $d_0 \neq 0$ is by assumption a descent direction for $-g$.

Employing the trust region $\Omega_j$ in the definition of $\beta_j$ guarantees the sufficient decrease condition

$$f(y_j) \leq f(x_k) - \frac{L}{2} \|x_k - y_j\|^2 \tag{4.3}$$

and monotonicity of the true objective $f$ during the SSC.

To see why (4.3) holds, notice that the second ball $\bar{B} = \bar{B}_{\|g\|/2L}(x_k + \frac{g}{2L})$ appearing in the definition of $\Omega_j$ does not depend on $j$, so that since $y_0 \in \bar{B}$ we have $y_j \in \bar{B}$ for every $j \in [0:T]$, with $T$ maximal iteration index of the SSC. This is enough to obtain (4.3) because for every $z \in \bar{B}$ we have

$$f(z) \leq f(\bar{x}) - \langle g, z - \bar{x} \rangle + \frac{L}{2} \|z - \bar{x}\|^2 \leq f(\bar{x}) - \frac{L}{2} \|\bar{x} - z\|^2\,, \tag{4.4}$$

where the first inequality is the standard descent lemma and the second follows from the definition of $\bar{B}$.

We prove that the true objective $f$ is monotone decreasing in the next lemma.

**Lemma 4.1.** *Let us assume $y_j \in \bar{B}_{\langle g, \hat{d}_j \rangle / L}(\bar{x})$. Then for every $\beta \in [0, \beta_j]$ we have*

$$\frac{d}{d\beta} f(y_j + \beta d_j) \leq 0\,,$$

*and thus in particular $f(y_j + \beta_j d_j) \leq f(y_j)$.*

*Proof.* We have

$$\frac{d}{d\beta} f(y_j + \beta d_j) = \|d_j\| \langle \nabla f(y_j + \beta d_j), \hat{d}_j \rangle$$

$$= \|d_j\| \langle (\nabla f(y_j + \beta d_j) + g) - g, \hat{d}_j \rangle = \|d_j\| (\langle \nabla f(y_j + \beta d_j) + g, \hat{d}_j \rangle - \langle g, \hat{d}_j \rangle)$$

$$\leq \|d_j\| (L \|\bar{x} - y_j - \beta d_j\| - \langle g, \hat{d}_j \rangle) \leq 0\,,$$

13

where we used $g = -\nabla f(\bar{x})$ and the Lipschitzianity of $\nabla f$ in the first inequality and

$$y_j + \beta_j d_j \in \bar{B}_{\langle g, \hat{d}_j \rangle / L}(\bar{x})$$

in the second. $\qquad\square$

The next result illustrates how the sequence $\{x_k\}$ generated by Algorithm 2 satisfies certain descent conditions. This is an adaptation to our setting of the ones used in the analysis of many proximal type gradient methods (see [4], [5], [13] and references therein). A subtle difference is the introduction of an "hidden sequence" $\{\tilde{x}_k\}$ to control the projection of the negative gradient on the tangent cone.

**Proposition 4.1.** *Let us consider the sequence $\{x_k\}$ generated by Algorithm 2 and assume that*

- *the angle condition* (3.2) *holds;*
- *the SSC condition terminates in a finite number of steps.*

*Then*

$$f(x_k) - f(x_{k+1}) \geq \frac{L}{2} \|x_k - x_{k+1}\|^2, \tag{4.5}$$

$$\|x_k - x_{k+1}\| \geq K \|\pi(T_\Omega(\tilde{x}_k), -\nabla f(\tilde{x}_k))\| \tag{4.6}$$

*for some $\tilde{x}_k \in \{y_j\}_{j=0}^T$ such that $f(x_{k+1}) \leq f(\tilde{x}_k) \leq f(x_k) - \frac{L}{2}\|x_k - \tilde{x}_k\|^2$, $\|\tilde{x}_k - x_k\| \leq \|x_{k+1} - x_k\|$ and for $K = \tau/(L(1+\tau))$.*

## 4.2  SSC for Frank-Wolfe variants

In this section, we show how to apply our results to the PFW, the AFW and the FDFW on polytopes, i.e., we prove finite termination of the SSC procedure when one of these methods is considered in Algorithm 2. We also give worst case and average worst case bounds for the number of iterations of the SSC. We start by proving a general termination criterion.

**Lemma 4.2.** *Assume that the method $\mathcal{A}$ applied to any linear function $L_g(x) = -\langle g, x \rangle$ on the feasible set $\Omega$ and with every stepsize maximal always terminates in at most $T$ iterations with an optimal solution, i.e. generates a sequence $\{y_j\}_{j \in [0,T']}$ with $T' \leq T$ and $y_{T'} \in \arg\min_{x \in \Omega} L_g(x)$. Then the SSC with the method $\mathcal{A}$ on the feasible set $\Omega$ always terminates in at most $T$ iterations.*

*Proof.* Assume by contradiction that the SSC does at least $T+1$ iterations, generating the sequence $\{y_j\}_{j \in [0:T+1]}$ before terminating. Notice that in this case the SSC must always do maximal steps for $j \in [0:T]$, because it terminates at step 8 when $\alpha_j = \beta_j$ and in particular if $\alpha_j < \alpha_{\max}^{(j)}$. Then for some $T' \leq T$ we must have that $y_{T'} \in \arg\min_{x \in \Omega} L_g(x)$, which gives a contradiction because in this case the method can't find a feasible descent direction in Phase I and terminates returning $y_{T'}$. $\qquad\square$

**Remark 2.** *Using the same line of reasoning, it is not difficult to prove that the SSC always terminates if the method $\mathcal{A}$ applied to linear objectives and with stepsizes always maximal generates a (possibly finite) sequence $\{y_j\}$ satisfying*

$$\liminf \pi_{y_j}(g) = 0. \tag{4.7}$$

14

We now denote with $\{S^{(j)}\}$ the sequence of active sets generated by the AFW and the PFW method in the SSC, and with $y_j$ proper convex combination of the elements in $S^{(j)}$. Furthermore, for the FDFW we assume that the maximal stepsize is given by feasibility conditions as in [26]:

$$\alpha_{\max}(x, d) = \max\{\alpha \in \mathbb{R}_{\geq 0} \mid x + \alpha d \in \Omega\}. \tag{4.8}$$

Notice that after a maximal in face step from $y_j$ we have $\dim(\mathcal{F}(y_{j+1})) < \dim(\mathcal{F}(y_j))$ because $y_{j+1}$ lies on the boundary of $\mathcal{F}(y_j)$.

**Proposition 4.2.** *The SSC always terminates in at most:*

- *$|A|$ iterations for the AFW,*
- *$|A| - 1$ iterations for the PFW,*
- *$\dim(\Omega) + 1$ iterations for the FDFW.*

*Proof.* By Lemma 4.2 we just need to bound the maximum number of iterations if the method performs always maximal steps for a linear objective $L_g(x)$. The AFW can do at most $|A| - 1$ consecutive maximal away steps, since at every such step the number of active atoms decreases by one. Analogously, the FDFW can do at most $\dim(\Omega)$ consecutive maximal in face steps, since at every such steps the dimension of the minimal face containing the current iterate decreases by one. The respective bound follows Lemma 4.2 by noticing that in the linear case the methods terminate after a full FW step. For the PFW, the linearity of the objective implies that only atoms in $\bar{A} := \operatorname{argmax}_{a \in A}\langle g, x \rangle$ can be added to the support, and only atoms in $A \setminus \bar{A}$ can be dropped from the support. In particular, once an atom is dropped from the active set it cannot be added again, and since at every maximal step the PFW drops an atom from the active set its maximal number of iterations is $|A \setminus \bar{A}| \leq |A| - 1$. $\qquad\square$

**Proposition 4.3.** *Assume that the linear minimizer is not changed during the SSC. Then, for an infinite sequence $\{x_k\}$, the worst case average number of iterations is*

- *2 for the AFW and the PFW,*
- *$\Delta(\Omega) + 1$ for the FDFW.*

The proof uses analogous arguments to the ones in [45, Theorem 8] to bound the number of bad steps and we defer it to the appendix.

## 4.3 Convergence rates

### 4.3.1 Smooth non convex objectives

We first prove, in the generic smooth non convex case, convergence to the set of stationary points with a rate of $O(\frac{1}{\sqrt{k}})$ for $\|\pi(T_\Omega(\tilde{x}_i), -\nabla f(\tilde{x}_i))\|$.

**Theorem 4.1.** *Let us consider the sequence $\{x_k\}$ generated by Algorithm 2 and assume that*

- *the angle condition (3.2) holds;*
- *the SSC procedure always terminates in a finite number of steps.*

15

| Algorithm | Article | LMO c.r. | Gradient c.r. | Gap |
|---|---|---|---|---|
| NCGS | [58] | $O\left(\frac{1}{k^{0.25}}\right)$ | $O\left(\frac{1}{\sqrt{k}}\right)$ | $\min_{0\le i\le k}\pi(x_i)$ |
| AFW, FW | [17], [44] | $O\left(\frac{1}{\sqrt{k}}\right)$ | $O\left(\frac{1}{\sqrt{k}}\right)$ | $\min_{0\le i\le k}G(x_i)$ |
| AFW, PFW, FDFW + SSC | Ours | $O\left(\frac{1}{\sqrt{k}}\right)$ | $O\left(\frac{1}{\sqrt{k}}\right)$ | $\min_{0\le i\le k}\|\pi(T_\Omega(\tilde{x}_i),-\nabla f(\tilde{x}_i))\|$ |

**Table 2:** Comparison between convergence rates in the generic smooth non convex case. See also Remark 3. $\pi(x)=\|x-\pi\left(\Omega,x-\frac{\nabla f(x)}{2L}\right)\|$, $G$ is the FW gap (see (4.10)).

*Then $\{f(x_k)\}$ is decreasing, $f(x_k)\to \tilde{f}\in\mathbb{R}$ and the limit points of $\{x_k\}$ are stationary. Furthermore, for any sequence $\{\tilde{x}_k\}$ satisfying the conditions of Proposition 4.1, we have $\|\tilde{x}_k-x_k\|\to 0$, and*

$$\min_{0\le i\le k}\|\pi(T_\Omega(\tilde{x}_i),-\nabla f(\tilde{x}_i))\|\le \min_{0\le i\le k}\frac{\|x_{i+1}-x_i\|}{K}\le \sqrt{\frac{2(f(x_0)-\tilde{f})}{K^2 L(k+1)}}, \qquad (4.9)$$

*for $K=\tau/(L(1+\tau))$.*

We now give a corollary for Theorem 4.1 specialized to the FW variants described in Section 3.1 (see also Table 2).

**Corollary 4.1.** *Let us assume that $\Omega=\mathrm{conv}(A)$, with $|A|<+\infty$ in Problem (2.1). Then the sequence $\{x_k\}$ generated by Algorithm 2 with AFW (PFW or FDFW) in the SSC converges at a rate given by equation (4.9), with $\tau=\tau_p/2$ ($\tau_p$ or $\tau_v/2$, respectively).*

*Proof.* Finite termination of the SSC follows by Proposition 4.2, and the angle condition is satisfied by Proposition 3.1. Thus we have all the assumptions to apply Theorem 4.1. ◻

**Remark 3.** *Let $G:\Omega\to\mathbb{R}_{\ge 0}$ be the FW gap (see, e.g., [44]):*

$$G(x)=\max_{s\in\Omega}\langle-\nabla f(x),s-x\rangle. \qquad (4.10)$$

*Then, for any $y\in\Omega$*

$$G(y)=\max_{s\in\Omega}\langle-\nabla f(y),s-y\rangle=\max_{s\in\Omega\setminus\{y\}}\|s-y\|\langle-\nabla f(y),\frac{s-y}{\|s-y\|}\rangle\le D\|\pi(T_\Omega(y),-\nabla f(y))\|,$$

(4.11)

*where the inequality follows from Proposition 2.2.*

Taking into account equation (4.11), it is easy to see that our rate is an improvement of the ones proved in [44] and [17] (see Table 2). Furthermore, we do not need to start from a vertex to avoid dependence from the support of $\{x_0\}$ like in [17, Theorem 5.1]. Finally, our method improves the conditional gradient sliding rate (NCGS) not only in LMO but also in gradients, given that from $\Omega-\{y\}\subset T_\Omega(y)$ it follows $\pi(y)\le \|\pi(T_\Omega(y),-\nabla f(y))\|/2L$ for every $y\in\Omega$.

### 4.3.2 Objectives with KL property

As a consequence of Proposition 4.1, we have linear convergence rates for the general algorithmic scheme reported in Algorithm 2 under the KL inequality (2.2), the angle condition (3.2), and finite termination of the SSC procedure. In the next results (Lemma 4.3, Theorem 4.2 and Corollary 4.2), we always assume the following:

- the angle condition (3.2) holds;
- the SSC procedure always terminates in a finite number of steps.

**Lemma 4.3.** *Let us consider the sequence $\{x_k\}$ generated by Algorithm 2 and assume that the objective function $f$ satisfies condition (2.2), with $f(x^*)$ fixed, in every feasible point generated by the algorithm. Then, for $q = \left(1 + \frac{\mu}{L}\frac{\tau^2}{(1+\tau)^2}\right)^{-1}$ we have $f(x_k) \to f(x^*)$, with*

$$f(x_k) - f(x^*) \le q^k (f(x_0) - f(x^*)), \tag{4.12}$$

*and $x_k \to \tilde{x}^*$ with*

$$\|x_k - \tilde{x}^*\| \le \frac{\sqrt{2-2q}(f(x_0) - f(\tilde{x}^*))}{\sqrt{L}(1 - \sqrt{q})} q^{\frac{k}{2}}, \tag{4.13}$$

*for $\tilde{x}^*$ stationary point such that $f(\tilde{x}^*) = f(x^*)$.*

As an example, the assumption of Lemma 4.3 is clearly satisfied if (2.2) holds globally, corresponding to a constrained version of the global PL property used in [40]. By [13, Corollary 6], for convex objectives this assumption is satisfied in particular under a global quadratic Holderian error bound, thus, e.g., by strongly convex objectives.
Under mild assumptions on the stationary point $x^*$, we can also apply Lemma 4.3 locally on non convex objectives, thus adapting to our projection free setting the local results given in [5, Section 2.3] for proximal methods.

**Theorem 4.2.** *Let Assumption 2.1 hold at $x^*$. Further assume that $x_k \in B_\delta(x^*) \Rightarrow f(x_{k+1}) \ge f(x^*)$. Then, for some $\tilde{\delta} > 0$, if $x_0 \in B_{\tilde{\delta}}(x^*)$ the rates (4.12) and (4.13) hold.*

It is not difficult to see that the assumption $x_k \in B_\delta(x^*) \Rightarrow f(x_{k+1}) \ge f(x^*)$ is true, e.g., if $x^*$ is a minimizer on its connected component of the sublevel set $[f \le f(x_0)]$.

As a corollary of Theorem 4.2, we can apply Lemma 4.3 and derive the following asymptotic rates.

**Corollary 4.2.** *Let us consider the sequence $\{x_k\}$ generated by Algorithm 2. Let Assumption 2.1 hold at every point of the limit set of $\{x_k\}$. Then, for some positive constants $M$ and $\bar{M}$, $\{x_k\} \to x^*$, with the asymptotic rates:*

$$\begin{aligned} f(x_k) - f(x^*) &\le M q^k, \\ \|x_k - x^*\| &\le \bar{M} q^{\frac{k}{2}}. \end{aligned} \tag{4.14}$$

Similarly to what we did for Theorem 4.1, here we give a corollary for Lemma 4.3 related to the FW variants described in Section 3.1.

**Corollary 4.3.** *Let us assume that the objective function $f$ satisfies condition (2.2) on every point generated by the algorithm, with $f(x^*)$ fixed, and that $\Omega = \text{conv}(A)$ with $|A| < +\infty$ in Problem (2.1). Then the sequence $\{x_k\}$ generated by Algorithm 2 with AFW (PFW or FDFW) in the SSC converges at the rates given by Lemma 4.3, with $\tau = \tau_p/2$ ($\tau_p$ or $\tau_v/2$, respectively).*

| Algorithm | Article | Objective | $\gamma(k)$ | $I_b$ | $q_{gs}$ | $h_k/h_0$ upper bound | $T_{avg}$ |
|---|---|---|---|---|---|---|---|
| AFW | [45] | SC | $k/2$ | $\lvert S_0\rvert-1$ | $1-\frac{\mu}{L}\frac{\tau_p^2}{4}$ | $\left(1-\frac{\mu}{L}\frac{\tau_p^2}{4}\right)^{\frac{k}{2}}$ | - |
| PFW | [45] | SC | $k/(3\lvert A\rvert!+1)$ | - | $1-\frac{\mu}{L}\tau_p^2$ | $\left(1-\frac{\mu}{L}\tau_p^2\right)^{\frac{k}{3\lvert A\rvert!+1}}$ | - |
| FDFW[1] | [43] | SC | $k/(\Delta(\Omega)+1)$ | $\dim(\mathcal{F}(x_0))$ | $1-\frac{\mu}{L}\frac{\tau_v^2}{4}$ | $\left(1-\frac{\mu}{L}\frac{\tau_v^2}{4}\right)^{\frac{k}{\Delta(\Omega)+1}}$ | - |
| AFW + SSC | Ours | NC, KL | $k$ | - | $\left(1+\frac{\mu}{L}\frac{\tau_p^2}{(2+\tau_p)^2}\right)^{-1}$ | $\left(1+\frac{\mu}{L}\frac{\tau_p^2}{(2+\tau_p)^2}\right)^{-k}$ | 2 |
| PFW + SSC | Ours | NC, KL | $k$ | - | $\left(1+\frac{\mu}{L}\frac{\tau_p^2}{(1+\tau_p)^2}\right)^{-1}$ | $\left(1+\frac{\mu}{L}\frac{\tau_p^2}{(1+\tau_p)^2}\right)^{-k}$ | 2 |
| FDFW + SSC | Ours | NC, KL | $k$ | - | $\left(1+\frac{\mu}{L}\frac{\tau_v^2}{(1+\tau_v)^2}\right)^{-1}$ | $\left(1+\frac{\mu}{L}\frac{\tau_v^2}{(1+\tau_v)^2}\right)^{-k}$ | $\Delta(\Omega)+1$ |

**Table 1:** Comparison between the rates of the standard and SSC version of some FW variants for $\Omega = \mathrm{conv}(A)$ with $\lvert A\rvert < \infty$. SC = strongly convex, NC = non convex, KL = KL property. $\gamma(k)$: lower bound on the number of good steps after $k$ steps, counting from the first good step. $I_b$: bound on the number of bad steps before the first good step. $q_{gs}$: rate in good steps. $h_k/h_0$ upper bound: worst case rate assuming no initial bad steps, equal to $q_{gs}^{\gamma(k)}$. $\Delta(\Omega)$ = maximum increase in face dimension $\mathcal{F}(x_{k+1}) - \mathcal{F}(x_k)$ after a FW step. $S_0$ = active set for $x_0$. $T_{avg}$ = worst case average iteration number of the SSC (see Proposition 4.3)

*Proof.* Finite termination of the SSC follows by Proposition 4.2, and the angle condition is satisfied by Proposition 3.1. Thus we have all the assumptions to apply Lemma 4.3. □

For comparison, we now recall some well-known result related to global linear convergence rates for the FW variants under analysis.

**Proposition 4.4.** *Let us assume that the objective function $f$ is $\mu-$strongly convex and $\Omega = \mathrm{conv}(A)$ with $\lvert A\rvert < +\infty$ in Problem (2.1). Let $\{x_k\}$ be a sequence generated by the AFW (PFW or FDFW), with stepsize given by exact linesearch. If the initial active set is $S_0 = \{x_0\}$ for the AFW ($S_0 = \{x_0\}$ for the PFW, $\dim(\mathcal{F}(x_0)) = 0$ for the FDFW), then*

$$f(x_k) - f^* \le q_{gs}^{\gamma(k)}(f(x_0) - f^*),  \tag{4.15}$$

*for $\gamma(k)$ and $q_{gs}$ given in Table 1.*

*Proof.* For the AFW and the PFW the result follows directly from [45, Theorem 1], with the exception of the good steps rate for the PFW, which can be obtained by applying the bound [45, Equation 10] in [45, Equation 5]. For the FDFW the result follows from [43, Theorem 1] (where the method is referred to as DiCG), with the bound $\mu\mathrm{PWidth}(V(\Omega)^2$ on the geometric strong convexity constant implied by [45, Theorem 6] improved to $\mu\mathrm{PFWidth}(\Omega)^2$ as in Proposition 3.1. □

For all the examples where an upper bound on $\tau_p = \frac{\mathrm{PWidth}(A)}{D}$ is known (see [59], [55] and references therein) when $\dim(\mathrm{conv}(A)) \to \infty$ then $\tau_p \to 0$ and our rates for the SSC converge to the rates without SSC for good steps in Table 1. While we are not able to prove this limit in general, for all polytopes with dimension greater or equal to 2, except low dimensional simplices (see Example 1), we still have $\tau_p \le \frac{1}{2}$

(because $\text{PdirW}(A,g,x) + \text{PdirW}(A,-g,x) \leq D$ for $x$ in the relative interior of $\text{conv}(A)$ and $\pm g$ feasible and orthogonal to $\text{conv}(S)$ for some $S \in S_x$). Using this together with Example 1 for simplices, it is easy to check that the rates in Corollary 4.3 (SSC based FW variants) are strict improvements on the known worst case rates (standard FW variants) reported in Proposition 4.4, with a limited number of exceptions. These are the trivial one dimensional case and simplices with low dimension ($\leq 4$ for the PFW, and $\leq 8$ for the AFW using the loose bounds in Example 1) combined with objectives having condition number $\mu/L$ sufficiently close to 1.

**Example 1.** *If $W(\text{conv}(A))$ is the width of $\text{conv}(A)$ (see [45, Section 3]) then it follows directly from the definition of $\text{PWidth}$ that $W(\text{conv}(A)) \geq \text{PWidth}(A)$, with equality for $A = \{e_1,...,e_n\}$ (see [45] and [55]). Let now $A = \{a_1,...,a_n\}$ be a set of $n$ affinely independent points in $\mathbb{R}^{n-1}$. We claim that, for $r_n = \sqrt{1 - \frac{1}{n}}$ circumradius of the $n-1$ dimensional unit simplex $\Delta_{n-1}$*

$$\text{PWidth}(A)/D \leq r_n^{-1} W(\Delta_{n-1}) = \begin{cases} 2r_n^{-1}\sqrt{\frac{1}{n}} & \text{for } n \text{ even}, \\ 2r_n^{-1}\sqrt{\frac{1}{n-1/n}} & \text{for } n \text{ odd}. \end{cases} \quad (4.16)$$

*To see this, assume without loss of generality $D = 1$ and $0 \in \text{int}(\Omega)$ for $\Omega = \text{conv}(A)$. Then if $\hat{A} = \{\hat{a}_1,...,\hat{a}_n\}$ we have $W(\text{conv}(\hat{A})) \geq W(\text{conv}(A))$. We can conclude*

$$\frac{\text{PWidth}(A)}{D} = \text{PWidth}(A) \leq W(\text{conv}(A)) \leq W(\text{conv}(\hat{A})) \leq r_n^{-1} W(\Delta_{n-1}), \quad (4.17)$$

*where in the last inequality we used that regular simplices maximize the width among simplices with fixed inradius (see, e.g., [3] and [30]).*

# 5 Examples

We now discuss some examples of objectives satisfying the KL property and sets where the angle condition can be satisfied with an explicit bound, relevant to practical optimization problems.

## 5.1 KL property

The KL property of Assumption 2.1 is satisfied for Problem (2.1) in the following cases:

- $f$ is composite strongly convex, i.e. $f(x) = g(Bx)$ with $g$ strongly convex, and $\Omega$ is a polytope [49, Proposition 4.1],

- $f$ is composite strongly convex as in the previous point, $\Omega$ is the $l^p$ ball for $p \in [1,2]$, and $\inf_{x \in \Omega} f(x) > \inf_{x \in \mathbb{R}^n} g(Bx)$ [49, Proposition 4.2],

- $f$ is (non convex) quadratic, i.e. $f(x) = x^\top Q x + b^\top x + c$, and $\Omega$ is a polytope, [49, Corollary 5.2],

- $f$ is non convex quadratic and does not satisfy the degeneracy condition of [34, equation (30)], and $\Omega$ is the unit sphere [34, Theorem 3.13].

## 5.2 Angle condition bounds

### 5.2.1 Bounds using PWidth

For the unit simplex and the unit cube explicit $\Theta(1/\sqrt{n})$ values were given in [55, Example 1 and 2]. With analogous arguments it can be proved that the PWidth of the $l_1$ ball is $1/\sqrt{n}$. By Proposition (3.1), this implies that the angle condition can be lower bounded with $\tau = \Theta(1/\sqrt{n})$ for the unit simplex and the $l_1$ ball, and with $\tau = \Theta(1/n)$ for the unit cube.

### 5.2.2 Bounds using facial distance vf

For a polytope $\Omega = \{x \in \mathbb{R}^n \mid Ax \leq b\}$ with $A \in \mathbb{R}^{m \times n}$ the facial distance can be defined as (see [8]):

$$\mathrm{vf}(\Omega) = \min_{\substack{v \in V(\Omega) \\ i : \langle a^{(i)}, v \rangle < b_i}} \frac{b_i - \langle a^{(i)}, v \rangle}{\|a^{(i)}\|}. \tag{5.1}$$

It is the easy to bound $\mathrm{vf}(\Omega)$ on some specific class of polytopes and, consequently, give an explicit bound for the angle condition (see also [7]). For instance, if the matrix $A$ is totally unimodular (i. e. all the vertices are integral for $b$ integral), we have the following properties.

**Proposition 5.1.** *If the matrix $A$ is totally unimodular and $b$ is integral, then for $\bar{a} = \max_{i \in [1:m]} \|a_i\|$:*

- *for the AFW or the PFW, if the size of the active set stays bounded by $\bar{s}$, then*

$$\mathrm{SB}_{\mathrm{AFW}}(\Omega) \geq \frac{1}{2\bar{s}\bar{a}D}, \quad \mathrm{SB}_{\mathrm{PFW}}(\Omega) \geq \frac{1}{\bar{s}\bar{a}D}; \tag{5.2}$$

- *for the FDFW,*

$$\mathrm{SB}_{\mathrm{FD}}(\Omega) \geq \frac{1}{2D\bar{a}(\dim(\Omega)+1)} \geq \frac{1}{2D\bar{a}(n+1)}. \tag{5.3}$$

*Proof.* If $A$ is totally unimodular then for $i \in [1:m], v \in V$ such that $b_i - \langle a^{(i)}, v \rangle > 0$ we have

$$\frac{b_i - \langle a^{(i)}, v \rangle}{\|a_i\|} \geq \frac{1}{\|a_i\|} \tag{5.4}$$

since the numerator on the LHS must be at least one. By applying (5.4) to the RHS of (5.1) we obtain

$$\mathrm{vf}(\Omega) \geq \min_{i \in [1:m]} \frac{1}{\|a_i\|} = \frac{1}{\bar{a}}. \tag{5.5}$$

Then the thesis follows for the AFW and the PFW directly from the bounds of Remark 1. For the FDFW, the second part of (5.3) is trivially true since $\dim(\Omega) \leq n$, and the first follows by the bound given in Remark 1, using that by the Caratheodory theorem for every feasible point $x$ there exists $S \in S_x$ with $|S| \leq \dim(\Omega) + 1$. $\square$

The bound of Proposition 5.1 allows us to bound the angle condition for the min cost flow polytope with integral capacities:

$$\Omega = \{x \in \mathbb{R}^n \mid Ax \leq b, \ 0 \leq x \leq c\}, \tag{5.6}$$

with $b, c$ integral and $A$ incidence matrix of a directed graph $G$.

**Corollary 5.1.** *Consider a directed graph $G$ with incidence matrix $A \in \mathbb{R}^{m \times n}$ and maximum degree of a vertex $d$. Then if $\Omega$ is given as in (5.6):*

$$\mathrm{SB}_{\mathrm{FD}}(\Omega) \geq \frac{1}{2\sqrt{d}(n+1)\|c\|} \tag{5.7}$$

*Proof.* By the capacity constraints, the diameter of $\Omega$ is at most $\|c\|$. Then the result follows easily from Proposition 5.1 by noticing that $\Omega$ can be rewritten as $\{x \in \mathbb{R}^n \mid \bar{A}x \leq b\}$ for $\bar{A} = (A; I; -I)$ totally unimodular (see, e.g., [63]) with maximum norm of a row equal to $\sqrt{d}$. $\qquad\square$

### 5.2.3 Bounds on sets with smooth boundary

On convex sets with smooth boundary the angle condition can be satisfied with constant arbitrarily close to 1 using orthographic retractions [60, Section 6.3]. Furthermore, on sublevel sets of smooth and strongly convex functions the FDFW satisfies the angle condition with constant equal to the condition number of the function divided by 2 [60, Section 6.2].

## 5.3 Applications

There is a number of practical optimization problems with the feasible sets and objectives discussed above. To start with, the LASSO problem, the minimum enclosing ball problem, training linear support vector machines and finding maximal cliques in graphs can all be formulated as convex quadratic optimization problems [18] on the $l_1$ ball or the simplex. The trust region subproblem is a non convex quadratic problem on the unit sphere (see [34]). The min cost flow problem with a quadratic objective is also of practical interest [62]. Many other examples can be found in [49].

# 6 Numerical tests

We tested the SSC on the AFW and the PFW methods, applied to a quadratic (non convex) relaxation of the maximum clique problem proposed in [15].
More precisely, let $A$ be the adjacency matrix of a graph $G$. In [15] it is proved that there is a one to one correspondence between the maximal cliques of $G$ and the local minima of the function $f : \Delta_{n-1} \to \mathbb{R}$ defined by

$$f(x) = -x^\mathsf{T} A x - \frac{1}{2}\|x\|^2. \tag{6.1}$$

Therefore, we consider instances of Problem (2.1) with objective (6.1) and feasible set the $n-1$ dimensional unit simplex, that is $\Omega = \Delta_{n-1}$.

**Table 3:** Max clique found, average clique size, standard deviation of clique sizes and average CPU time for AFW and SSC + AFW on max clique instances from the DIMACS benchmark.

| Instance | AFW | | | | SSC + AFW | | | |
|---|---|---|---|---|---|---|---|---|
| | Max | Mean | Std | CPU time | Max | Mean | Std | CPU time |
| C2000.5 | 14 | 11.7 | 0.89 | 2.800 | 14 | 11.6 | 1.00 | 0.082 |
| C2000.9 | 67 | 60.2 | 2.20 | 3.135 | 65 | 60.0 | 2.05 | 0.200 |
| C4000.5 | 16 | 12.8 | 0.94 | 23.487 | 16 | 12.5 | 0.92 | 0.429 |
| MANN_a81 | 1080 | 1080.0 | 0.00 | 31.156 | 1080 | 1080.0 | 0.00 | 25.047 |
| keller6 | 45 | 38.4 | 2.41 | 13.713 | 43 | 37.8 | 2.22 | 0.413 |

**Table 4:** Max clique found, average clique size, standard deviation of clique sizes and average CPU time for PFW and SSC + PFW on max clique instances from the DIMACS benchmark.

| Instance | PFW | | | | SSC + PFW | | | |
|---|---|---|---|---|---|---|---|---|
| | Max | Mean | Std | CPU time | Max | Mean | Std | CPU time |
| C2000.5 | 14 | 11.8 | 0.86 | 2.811 | 14 | 12.1 | 0.86 | 0.077 |
| C2000.9 | 67 | 62.3 | 1.83 | 3.031 | 68 | 62.0 | 1.77 | 0.150 |
| C4000.5 | 15 | 12.7 | 0.92 | 23.423 | 16 | 13.4 | 0.95 | 0.379 |
| MANN_a81 | 1080 | 1080.0 | 0.00 | 19.867 | 1080 | 1080.0 | 0.00 | 15.442 |
| keller6 | 44 | 37.3 | 2.68 | 13.515 | 45 | 35.6 | 2.83 | 0.258 |

The graph instances we use are taken from the DIMACS benchmark [36]. To have a fair comparison for both the AFW and the PFW we use the stepsize given by

$$\alpha_k = \min\{\alpha_k^{\max}, -\frac{\langle \nabla f(x_k), d_k \rangle}{L\|d_k\|^2}\} \tag{6.2}$$

with $\alpha_k^{\max}$ determined by boundary conditions. In this way the new point computed by the methods coincides with the first point computed in the SSC procedure of their multistep versions.

We reported in Table 3, 4 the results for the most challenging instances, aggregated on 100 runs starting from random points. The SSC clearly improves the CPU times while keeping the solution quality. Indeed in these problems the SSC allows the methods to identify the support of a local minimum in fewer iterations, so that the slow initial convergence phase is skipped (see Figures 1, 2).
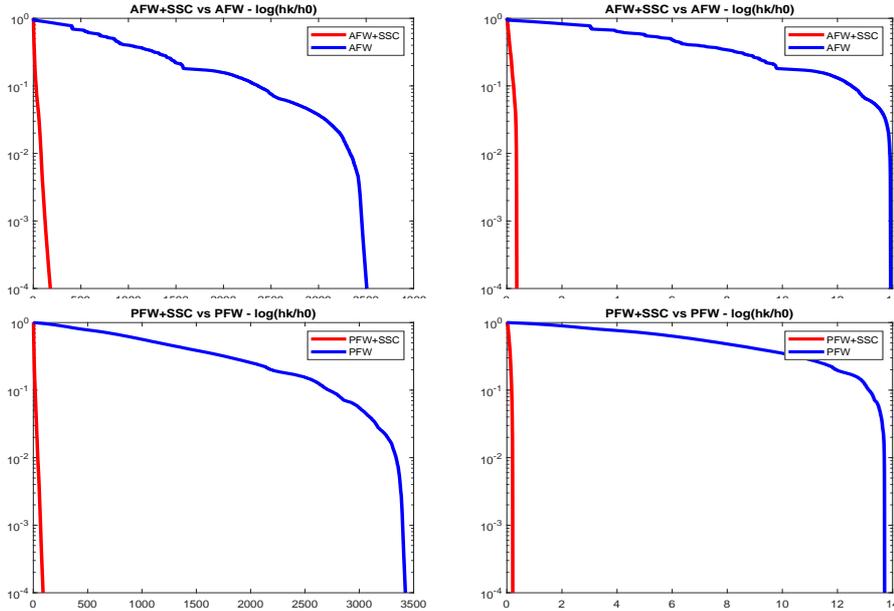
Figure 1: Iteration number and CPU time vs $\log(h_k/h_0)$ in the first and the second column respectively for the instance keller6

## 7 Conclusions

FW variants rely on the choice of good feasible descent directions, for which there needs to be a trade-off between slope and maximal stepsize. To address this issue we proposed the SSC procedure, which allowed us to prove bad step free convergence rates under an angle condition for the directions selected by the method. Preliminary numerical experiments also support the soundness of this approach.

Future research directions include employing our framework to design and analyze other projection free first order methods, investigating active set identification properties of FW variants with the SSC, generalizing our framework to constrained stochastic optimization, as well as applications for the solution of real-world data science problems.

## 8 Appendix

### 8.1 KL property

We state here a result showing an implication between the (global) PL property used in [40] and (2.2). We first recall the PL property used in [40]:

$$\frac{1}{2}\|\nabla f(x)\|^2 \geq \mu(f(x) - f^*).\qquad(8.1)$$

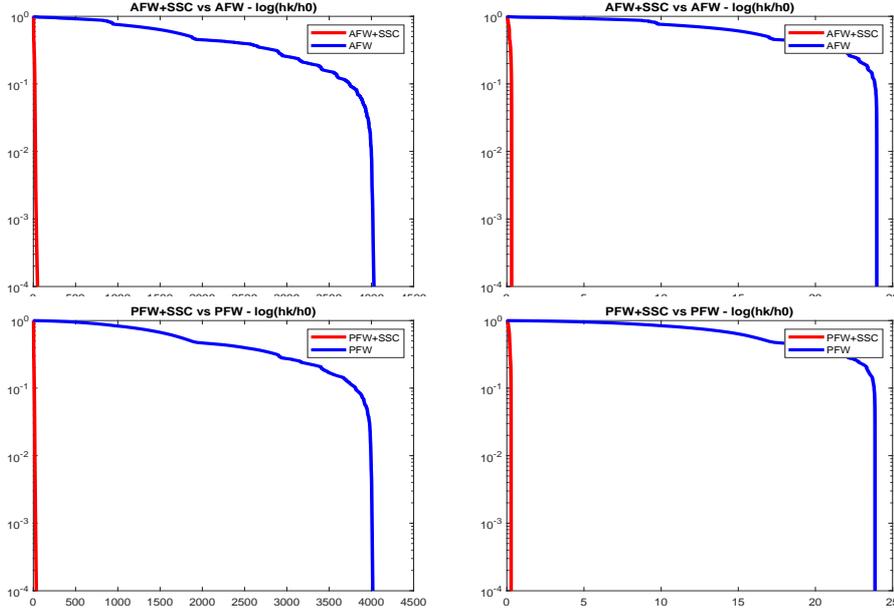with $f^*$ optimal value of $f$ with non empty solution set $\mathcal{X}^*$.

Figure 2: Iteration number and CPU time vs $\log(h_k/h_0)$ in the first and the second column respectively for the instance C4000.5

**Proposition 8.1.** *If $f$ is convex, the optimal solution set $\mathcal{X}^*$ of $f$ is contained in $\Omega$ and (8.1) holds, then (2.2) holds for every $x \in \Omega$.*

*Proof.* By [40, Theorem 2] the PL property is equivalent, for convex objectives, to the unconstrained quadratic growth condition:

$$f(x) - f^* \geq \frac{\mu}{2} \mathrm{dist}(x, \mathcal{X}^*)^2 \tag{8.2}$$

In turn, given that by the assumption $\mathcal{X}^* \subset \Omega$ the set $\mathcal{X}^*$ is the solution set for $f_\Omega$ as well, (8.2) implies the global non smooth Holderian error bound condition from [13] with $\varphi(t) = \sqrt{\frac{2t}{\mu}}$, and by [13, Corollary 6] this is equivalent to the KL property (2.2) holding globally on $\Omega$. $\qquad\square$

**Remark 4.** *We remark that without the assumption $\mathcal{X}^* \subset \Omega$ the implication is no longer true even for convex objectives, a counter example being $\Omega$ equal to the unitary ball and $f((x^{(1)}, ..., x^{(n)})) = (x^{(1)} - 1)^2$. At the same time, the KL property we used does not imply the PL property in general, since the latter only deals with unconstrained minima.*

## 8.2 Proofs

We report here the missing proofs. We start with the proof of Lemma 3.2.

*Proof.* By the standard descent lemma [10, Proposition 6.1.2],

$$f(x_{k+1}) = f(x_k + \alpha_k d_k) \leq f(x_k) + \alpha_k \langle \nabla f(x_k), d_k \rangle + \alpha_k^2 \frac{L}{2} \|d_k\|^2, \tag{8.3}$$

24

and in particular

$$f(x_k)-f(x_{k+1})\geq -\alpha_k\langle\nabla f(x_k),d_k\rangle-\alpha_k^2\frac{L}{2}\|d_k\|^2\geq\frac{L}{2}\alpha_k^2\|d_k\|^2=\frac{L}{2}\|x_{k+1}-x_k\|^2\,,\quad(8.4)$$

where we used $\alpha_k\leq\bar{\alpha}_k$ in the last inequality. This proves (3.13). $\qquad\square$

We now state a preliminary result needed to prove Proposition 2.2:

**Proposition 8.2.** *Let $C$ be a closed convex cone. For every $y\in\mathbb{R}^n$*

$$\mathrm{dist}(C^*,y)=\sup_{c\in C}\langle\hat{c},y\rangle\,.$$

As stated in [21] this is an immediate consequence of the Moreau-Yosida decomposition:

$$y=\pi(C,y)+\pi(C^*,y)\,.$$

*Proposition 2.2.* First, by continuity of the scalar product we have

$$\sup_{h\in\Omega/\{\bar{x}\}}\left(g,\frac{h-\bar{x}}{\|h-\bar{x}\|}\right)=\sup_{h\in T_\Omega(\bar{x})\setminus\{0\}}(g,\hat{h})\,.\quad(8.5)$$

Since $N_\Omega(\bar{x})=T_\Omega(\bar{x})^*$ the first equality is exactly the one of Proposition 8.2 if $g\notin N_\Omega(\bar{x})$, and it is trivial since both terms are clearly 0 if $g\in N_\Omega(\bar{x})$.
It remains to prove

$$\mathrm{dist}(N_\Omega(\bar{x}),g)=\|\pi(T_\Omega(\bar{x}),g)\|\,,$$

which is true by the Moreau - Yosida decomposition. $\qquad\square$

*Proposition 4.1.* Let $B_j=\bar{B}_{\langle g,\hat{d}_j\rangle/L}(x_k)$ and let $T$ be such that $x_{k+1}=y_T$.
Inequality (4.3) applied with $j=T$ gives (4.5). Moreover, by taking $\tilde{x}_k=y_{\tilde{T}}$ for some $\tilde{T}\in[0:T]$ the conditions

$$f(x_{k+1})\leq f(\tilde{x}_k)\leq f(x_k)-\frac{L}{2}\|x_k-\tilde{x}_k\|^2\quad(8.6)$$

are satisfied by Lemma 4.1 and (4.3).
Let now $p_j=\|\pi(T_\Omega(y_j),-\nabla f(y_j))\|$ and $\tilde{p}_j=\|\pi(T_\Omega(y_j),g)\|=\|\pi(T_\Omega(y_j),-\nabla f(x_k))\|$.
We have

$$|p_j-\tilde{p}_j|\leq L\|y_j-x_k\|\,,\quad(8.7)$$

reasoning as for (3.17). We now distinguish four cases according to how the SSC terminates.
**Case 1:** $T=0$ or $d_T=0$. Since there are no descent directions $x_{k+1}=y_T$ must be stationary for the gradient $g$. Equivalently, $\tilde{p}_T=\|\pi(T_\Omega(x_{k+1}),g)\|=0$. We can now write

$$\|x_{k+1}-x_k\|\geq\frac{1}{L}(|p_T-\tilde{p}_T|)=\frac{p_T}{L}>Kp_T\,,$$

where we used (8.7) in the first inequality and $\tilde{p}_T=0$ in the equality. Finally, it is clear that if $T=0$ then $d_0=0$, since $y_0$ must be stationary for $-g$.
Before examining the remaining cases we remark that if the SSC terminates in Phase II then $\alpha_{T-1}=\beta_{T-1}$ must be maximal w.r.t. the conditions $y_T\in B_{T-1}$ or $y_T\in\bar{B}$. If

25

$\alpha_{T-1} = 0$ then $y_{T-1} = y_T$, and in this case we cannot have $y_{T-1} \in \partial \bar{B}$, otherwise the SSC would terminate in Phase II of the previous cycle. Therefore necessarily $y_T = y_{T-1} \in \text{int}(B_{T-1})^c$ (Case 2). If $\beta_{T-1} = \alpha_{T-1} > 0$ we must have $y_{T-1} \in \Omega_{T-1} = B_{T-1} \cap \bar{B}$, and $y_T \in \partial B_{T-1}$ (case 3) or $y_T \in \partial \bar{B}$ (case 4) respectively.

**Case 2:** $y_{T-1} = y_T \in \text{int}(B_{T-1})^c$. We can rewrite the condition as

$$\langle g, \hat{d}_{T-1} \rangle \leq L\|y_{T-1} - x_k\| = L\|y_T - x_k\|. \tag{8.8}$$

Thus

$$p_T = p_{T-1} \leq \tilde{p}_{T-1} + L\|y_T - x_k\| \leq \frac{1}{\tau}\langle g, \hat{d}_{T-1} \rangle + L\|y_T - x_k\| \leq \left(\frac{L}{\tau} + L\right)\|y_T - x_k\|, \tag{8.9}$$

where in the equality we used $y_T = y_{T-1}$, the first inequality follows from (8.7) and again $y_T = y_{T-1}$, the second from $\frac{\langle g, \hat{d}_T \rangle}{\tilde{p}_T} \geq \text{DSB}_{\mathcal{A}}(\Omega, y_T, g) \geq \text{SB}_{\mathcal{A}}(\Omega) = \tau$, and the third from (8.8). Then $\tilde{x}_k = x_{k+1} = y_T$ satisfies the desired conditions.

**Case 3:** $y_T = y_{T-1} + \beta_{T-1} d_{T-1}$ and $y_T \in \partial B_{T-1}$. Then from $y_{T-1} \in B_{T-1}$ it follows

$$L\|y_{T-1} - x_k\| \leq \langle g, \hat{d}_{T-1} \rangle, \tag{8.10}$$

and $y_T \in \partial B_{T-1}$ implies

$$\langle g, \hat{d}_{T-1} \rangle = L\|y_T - x_k\|. \tag{8.11}$$

Combining (8.10) with (8.11) we obtain

$$L\|y_{T-1} - x_k\| \leq L\|y_T - x_k\|. \tag{8.12}$$

Thus

$$p_{T-1} \leq \tilde{p}_{T-1} + L\|y_{T-1} - x_k\| \leq \frac{1}{\tau}\langle g, \hat{d}_{T-1} \rangle + L\|y_{T-1} - x_k\| \leq \left(\frac{L}{\tau} + L\right)\|y_T - x_k\|,$$

where we used (8.11), (8.12) in the last inequality and the rest follows reasoning as for (8.9). In particular we can take $\tilde{x}_k = y_{T-1}$, where $\|\tilde{x}_k - x_k\| \leq \|x_{k+1} - x_k\|$ by (8.12).

**Case 4:** $y_T = y_{T-1} + \beta_{T-1} d_{T-1}$ and $y_T \in \partial \bar{B}$.

The condition $x_{k+1} = y_T \in \bar{B}$ can be rewritten as

$$L\|x_{k+1} - x_k\|^2 - \langle g, x_{k+1} - x_k \rangle = 0. \tag{8.13}$$

For every $j \in [0:T]$ we have

$$x_{k+1} = y_j + \sum_{i=j}^{T-1} \alpha_i d_i. \tag{8.14}$$

We now want to prove that for every $j \in [0:T]$

$$\|x_{k+1} - x_k\| \geq \|y_j - x_k\|. \tag{8.15}$$

Indeed, we have

$$L\|x_{k+1} - x_k\|^2 = \langle g, x_{k+1} - x_k \rangle = \langle g, y_j - x_k \rangle + \sum_{i=j}^{T-1} \alpha_i \langle g, d_i \rangle$$

$$\geq \langle g, y_j - x_k \rangle \geq L\|y_j - x_k\|^2,$$

26

where we used (8.13) in the first equality, (8.14) in the second, $\langle g, d_j \rangle \geq 0$ for every $j$ in the first inequality and $y_j \in \bar{B}$ in the second inequality.

We also have

$$\frac{\langle g, x_{k+1} - x_k \rangle}{\|x_{k+1} - x_k\|} = \frac{\langle g, \sum_{j=0}^{T-1} \alpha_j d_j \rangle}{\| \sum_{j=0}^{T-1} \alpha_j d_j \|} \geq \frac{\langle g, \sum_{j=0}^{T-1} \alpha_j d_j \rangle}{\sum_{j=0}^{T-1} \alpha_j \|d_j\|}$$

$$\geq \min \left\{ \frac{\langle g, d_j \rangle}{\|d_j\|} \mid 0 \leq j \leq T-1 \right\}. \tag{8.16}$$

Thus for $\tilde{T} \in \operatorname{argmin} \left\{ \frac{\langle g, d_j \rangle}{\|d_j\|} \mid 0 \leq j \leq T-1 \right\}$

$$\langle g, \hat{d}_{\tilde{T}} \rangle \leq \frac{\langle g, x_{k+1} - x_k \rangle}{\|x_{k+1} - x_k\|} = L\|x_{k+1} - x_k\|, \tag{8.17}$$

where we used (8.16) in the first inequality and (8.13) in the second.

We finally have

$$p_{\tilde{T}} \leq \tilde{p}_{\tilde{T}} + L\|y_{\tilde{T}} - x_k\| \leq \frac{1}{\tau} \langle g, \hat{d}_{\tilde{T}} \rangle + L\|y_{\tilde{T}} - x_k\| \leq \left( \frac{L}{\tau} + L \right) \|x_{k+1} - x_k\|,$$

where we used (8.15), (8.17) in the last inequality and the rest follows reasoning as for (8.9). In particular $\tilde{x}_k = y_{\tilde{T}}$ satisfies the desired properties, where $\|\tilde{x}_k - x_k\| \leq \|x_{k+1} - x_k\|$ by (8.15). $\qquad \square$

*Proof of Proposition 4.3.* Let $T(k)$ be the number of iterates generated by the SSC at the step $k$ in Phase II. For the AFW and the PFW, reasoning as in the proof of Proposition 4.2 we obtain that if the SSC does $T(k)$ iterations, the number of active vertices decreases by at least $T(k) - 2$. Then on the one hand

$$|S^{(k)}| - |S^{(0)}| \geq 1 - |S^{(0)}|, \tag{8.18}$$

while on the other hand

$$|S^{(k)}| - |S^{(0)}| = \sum_{i=0}^{k-1} (|S^{(i+1)}| - |S^{(i)}|)$$

$$\leq 2k - \sum_{i=0}^{k-1} T(i). \tag{8.19}$$

Combining (8.18) and (8.19) and rearranging, we obtain:

$$\frac{1}{k} \sum_{i=0}^{k-1} T(i) \leq 2 + \frac{|S^{(0)}| - 1}{k}, \tag{8.20}$$

and the desired result follows by taking the limit for $k \to \infty$.

For the FDFW, notice that at every iteration the SSC performs a sequence of maximal in face steps terminated either by a Frank Wolfe step, after which $\mathcal{F}(y_j)$ can increase of at most $\Delta(\Omega)$, or by a non maximal in face step, after which $\mathcal{F}(y_j)$ stays the same. In both cases, we have

$$\dim(\mathcal{F}(x_{k+1})) - \dim(\mathcal{F}(x_k)) \leq \Delta(\Omega) - T(k) + 1. \tag{8.21}$$

27

Then,
$$\dim \mathcal{F}(x_k) - \dim \mathcal{F}(x_0) \geq -\dim \mathcal{F}(x_0), \tag{8.22}$$

and
$$\dim \mathcal{F}(x_k) - \dim \mathcal{F}(x_0) = \sum_{i=0}^{k-1} (\dim(\mathcal{F}(x_{i+1}) - \dim(\mathcal{F}(x_i))))$$
$$\leq k\Delta(\Omega) + k - \sum_{i=0}^{k-1} T(i). \tag{8.23}$$

The conclusion follows as for the AFW and the PFW. $\qquad\square$

*Theorem 4.1.* The sequence $\{f(x_k)\}$ is decreasing by (4.5). Thus by compactness $f(x_k) \to \tilde{f} \in \mathbb{R}$ and in particular $f(x_k) - f(x_{k+1}) \to 0$. So that by (4.5) also $\|x_{k+1} - x_k\| \to 0$. Let $\{x_{k(i)}\} \to \tilde{x}^*$ be any convergent subsequence of $\{x_k\}$. For $\{\tilde{x}_k\}$ chosen as in the proof of Proposition 4.1 we have $\|\tilde{x}_k - x_k\| \leq \|x_{k+1} - x_k\|$ because $\tilde{x}_k = y_T = x_k$ in case 1 and case 2, by (8.12) in case 3, and by (8.15) in case 4. Therefore

$$\|\tilde{x}_{k(i)} - x_{k(i)}\| \leq \|x_{k(i)+1} - x_{k(i)}\| \to 0.$$

Furthermore, $\|\pi(T_\Omega(\tilde{x}_{k(i)}), -\nabla f(\tilde{x}_{k(i)})))\| \leq \frac{\|x_{k(i)+1} - x_{k(i)}\|}{K} \to 0$ again by Proposition 4.1, so that $\tilde{x}_{k(i)} \to \tilde{x}^*$ with $\|\pi(T_\Omega(\tilde{x}_{k(i)}), -\nabla f(\tilde{x}_{k(i)}))\| \to 0$. Then $\|\pi(T_\Omega(\tilde{x}^*), -\nabla f(\tilde{x}^*))\| = 0$ and $\tilde{x}^*$ is stationary.

The first inequality in (4.9) follows directly from (4.6). As for the second, we have

$$\frac{k+1}{K^2} (\min_{0 \leq i \leq k} \|x_{i+1} - x_i\|)^2 = \frac{k+1}{K^2} \min_{0 \leq i \leq k} \|x_{i+1} - x_i\|^2$$
$$\leq \frac{1}{K^2} \sum_{i=0}^{k} \|x_i - x_{i+1}\|^2 \leq \frac{2}{LK^2} \sum_{i=0}^{k} (f(x_{i+1}) - f(x_i)) \leq \frac{2(f(x_0) - \tilde{f})}{LK^2},$$

where we used (4.5) in the first inequality, $\{f(x_i)\}$ decreasing together with $f(x_i) \to \tilde{f}$ in the second and the thesis follows by rearranging terms. $\qquad\square$

We now prove Lemma 4.3. We start by recalling Karamata's inequality ([38], [39]) for concave functions. Given $A, B \in \mathbb{R}^N$ it is said that $A$ majorizes $B$, written $A \succ B$, if

$$\sum_{i=1}^{j} A_i \geq \sum_{i=1}^{j} B_i \text{ for } j \in [1 : N],$$
$$\sum_{i=1}^{N} A_i = \sum_{i=1}^{N} B_i.$$

If $h$ is concave and $A \succ B$ by Karamata's inequality

$$\sum_{i=1}^{N} h(A_i) \leq \sum_{i=1}^{N} h(B_i).$$

In order to prove Lemma 4.3 we first need the following technical Lemma.

28

**Lemma 8.1.** Let $\{\tilde{f}_i\}_{i \in [0:j]}$ be a sequence of nonnegative numbers such that $\tilde{f}_{i+1} \le q\tilde{f}_i$ for some $q < 1$. Then

$$\sum_{i=0}^{j-1} \sqrt{\tilde{f}_i - \tilde{f}_{i+1}} \le \frac{\sqrt{\tilde{f}_0(1-q)}}{1-\sqrt{q}}. \tag{8.24}$$

*Proof.* Let $\bar{j} = \max\{i \ge 0 \mid \tilde{f}_j \le q^i\tilde{f}_0\}$, so that by (8.32) we have $\bar{j} \ge j$. Define $w^*, v \in \mathbb{R}_{\ge 0}^{\bar{j}+1}$ by

$$\begin{aligned}
v &= (\tilde{f}_0 - q\tilde{f}_0, ..., q^{\bar{j}-1}\tilde{f}_0 - q^{\bar{j}}\tilde{f}_0, q^{\bar{j}}\tilde{f}_0 - \tilde{f}_j), \\
w^* &= (\tilde{f}_0 - \tilde{f}_1, ..., \tilde{f}_{j-1} - \tilde{f}_j, 0, ..., 0).
\end{aligned} \tag{8.25}$$

Then for $0 \le l < \bar{j}$ we have

$$\sum_{i=0}^{l} v_i = \tilde{f}_0 - q^{l+1}\tilde{f}_0 \le \tilde{f}_0 - \tilde{f}_{\min(l+1,j)} = \sum_{i=0}^{l} w_i^*, \tag{8.26}$$

where we used $q^{l+1}\tilde{f}_0 \ge \tilde{f}_{l+1}$ for $l \le j-1$ and $q^{l+1}\tilde{f}_0 \ge \tilde{f}_j$ for $j \le l < \bar{j}$ in the inequality. Furthermore, for $l = \bar{j}$ we have

$$\sum_{i=0}^{l} v_i = \tilde{f}_0 - \tilde{f}_j = \sum_{i=0}^{l} w_i^*. \tag{8.27}$$

Now if $w$ is the permutation in descreasing order of $w^*$, clearly thanks to (8.26), and (8.27) we have $w \succ v$. Then

$$\begin{aligned}
\sum_{i=0}^{j-1} \sqrt{\tilde{f}_i - \tilde{f}_{i+1}} = \sum_{i=0}^{\bar{j}+1} \sqrt{w_i^*} &= \sum_{i=0}^{\bar{j}+1} \sqrt{w_i} \le \sum_{i=0}^{\bar{j}+1} \sqrt{v_i} \\
&\le \sqrt{\tilde{f}_0} \sum_{i=0}^{+\infty} \sqrt{q^i - q^{i+1}} = \frac{\sqrt{\tilde{f}_0(1-q)}}{1-\sqrt{q}},
\end{aligned} \tag{8.28}$$

where the first inequality follows from Karamata's inequality. $\square$

*Proof of Lemma 4.3.* If the sequence $\{x_k\}$ is finite, with $x_m = \tilde{x}$ stationary for some $m \ge 0$, we define $x_k = x_m$ for every $k \ge m$, so that we can always assume $\{x_k\}$ infinite. Notice that with this convention the sufficient decrease condition (4.5) is still satisfied for every $k$. Let $f_k = f(x_k) - f(x^*)$. $\{f_k\}$ is monotone decreasing by (4.5), and nonnegative since (2.2) holds for every $x_k$.

We want prove $f_{k+1} \le qf_k$. This is clear if $f_{k+1} = 0$. Otherwise using the notation of Proposition 4.1 we have

$$f_k - f_{k+1} \ge \frac{L}{2}\|x_k - x_{k+1}\|^2 \ge \frac{LK^2}{2}\|\pi(T_\Omega(\tilde{x}_k), -\nabla f(\tilde{x}_k))\|, \tag{8.29}$$

where we used (4.5) in the first inequality, (4.6) in the second. Since $\tilde{x}_k \in \{y_j\}_{j=0}^T$ by Proposition 4.1, we can apply (2.2) in $\tilde{x}_k$ to obtain

$$\frac{LK^2}{2}\|\pi(T_\Omega(\tilde{x}_k), -\nabla f(\tilde{x}_k))\|^2 \ge \mu LK^2(f(\tilde{x}_k) - f(x^*)) \ge \mu LK^2 f_{k+1}. \tag{8.30}$$

29

Concatenating (8.29), (8.30) and rearranging we obtain

$$f_{k+1} \leq (1 + \mu LK^2)^{-1} f_k = q f_k. \tag{8.31}$$

Thus by induction for any $i \geq 0$

$$f_{k+i} \leq q^i f_k, \tag{8.32}$$

which implies in particular (4.12).
We can now bound the length of the tails of $\{x_k\}$:

$$\sum_{i=0}^{+\infty} \|x_{k+i} - x_{k+i+1}\| \leq \sqrt{\frac{2}{L}} \sum_{i=0}^{+\infty} \sqrt{f_{k+i} - f_{k+i+1}} \leq \frac{\sqrt{2 f_k (1-q)}}{\sqrt{L}(1 - \sqrt{q})} \leq \frac{\sqrt{2 f_0 (1-q)}}{\sqrt{L}(1 - \sqrt{q})} q^{\frac{k}{2}}, \tag{8.33}$$

where we used (4.5) in the first inequality, Lemma 8.1 with $\{\tilde{f}_i\} = \{f_{k+i}\}$ and for $j \to +\infty$ in the second inequality, and (8.32) in the third. In particular $x_k \to \tilde{x}^*$ with

$$\|x_k - \tilde{x}^*\| \leq \sum_{j=0}^{+\infty} \|x_{k+j} - x_{k+j+1}\| = \frac{\sqrt{2 f_0 (1-q)}}{\sqrt{L}(1 - \sqrt{q})} q^{\frac{k}{2}} \tag{8.34}$$

by (8.33). $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \square$

*Proof of Theorem 4.2.* By continuity, for $\tilde{\delta} \to 0$ and $f_0 = f(x_0) - f(x^*)$ we have that

$$\max_{x_0 \in B_{\tilde{\delta}}(x^*) \cap [f \geq f(x^*)]} f_0 \to 0, \tag{8.35}$$

so we can take $\tilde{\delta} < \delta/2$ small enough in such a way that

$$\max_{x_0 \in B_{\tilde{\delta}}(x^*) \cap [f \geq f(x^*)]} \frac{\sqrt{2 f_0 (1-q)}}{L(1 - \sqrt{q})} + \sqrt{\frac{2}{L}} \sqrt{f_0} < \frac{\delta}{2}. \tag{8.36}$$

Let now $x_0 \in B_{\tilde{\delta}}(x^*) \cap [f \geq f(x^*)]$, so that

$$\tilde{\delta} < \frac{\delta}{2} < \delta - \frac{\sqrt{2 f_0 (1-q)}}{L(1 - \sqrt{q})} - \sqrt{\frac{2}{L}} \sqrt{f_0}, \tag{8.37}$$

where we use (8.36) in the second inequality. We now want to prove, by induction on $k$, $\{x_i\}_{i \in [0:k]} \subset B_\delta(x^*)$ with $f(x_{i+1}) \leq q f(x_i)$ for every $i \in [0:k]$ and $k \in \mathbb{N}$. To start with,

$$\sum_{i=0}^{k-1} \|x_i - x_{i+1}\| \leq \sqrt{\frac{2}{L}} \sum_{i=0}^{k-1} \sqrt{f_i - f_{i+1}} \leq \frac{\sqrt{2 f_0 (1-q)}}{\sqrt{L}(1 - \sqrt{q})} \tag{8.38}$$

where we used (4.5) in the first inequality, and Lemma 8.1 (which we can apply thanks to the inductive assumption) in the second. But then

$$\|x_{k+1} - x^*\| \leq \|x_0 - x^*\| + \left( \sum_{i=0}^{k-1} \|x_i - x_{i+1}\| \right) + \|x_k - x_{k+1}\|$$

$$\leq \tilde{\delta} + \frac{\sqrt{2 f_0 (1-q)}}{L(1 - \sqrt{q})} + \sqrt{\frac{2}{L}} \sqrt{f_k - f_{k+1}} \tag{8.39}$$

$$< \tilde{\delta} + \frac{\sqrt{2 f_0 (1-q)}}{L(1 - \sqrt{q})} + \sqrt{\frac{2}{L}} \sqrt{f_k} < \delta,$$

30

where we used (8.38) together with (4.5) in the second inequality, the assumption $x_k \in B_\delta(x^*) \Rightarrow f_{k+1} \geq 0$ in the third inequality, and (8.37) together with $f_0 \geq f_k$ in the last inequality.

We now have

$$
\begin{aligned}
\|\tilde{x}_k - x^*\| &\leq \|x_0 - x^*\| + \left(\sum_{i=0}^{k-1} \|x_i - x_{i+1}\|\right) + \|x_k - \tilde{x}_k\| \\
&\leq \|x_0 - x^*\| + \left(\sum_{i=0}^{k-1} \|x_i - x_{i+1}\|\right) + \|x_k - x_{k+1}\| < \delta,
\end{aligned}
\tag{8.40}
$$

where we use $\|\tilde{x}_k - x_k\| \leq \|x_{k+1} - x_k\|$ in the second inequality and the last inequality follows as in (8.40). Thus $\tilde{x}_k \in B_\delta(x^*)$ as well, which is enough to prove (8.31) and complete the induction. We have thus obtained $\{\tilde{x}_k\}, \{x_k\} \subset B_\delta(x^*)$, and the conclusion follows exactly as in the proof of Lemma 4.3. □

*Proof of Corollary 4.2.* Let $x^*$ be a limit point of $\{x_k\}$, and let $\tilde{\delta}$ be as in Theorem 4.2. First, for some $\bar{k} \in \mathbb{N}$ we must have $x_{\bar{k}} \in B_{\tilde{\delta}}(x^*)$. Furthermore, for every $k \in \mathbb{N}$ we have $f(x_k) \geq f(x^*)$ because $f(x_k)$ is non increasing and converges to $f(x^*)$. Thus we have all the necessary assumptions to obtain the asymptotic rates by applying Theorem 4.2 to $\{y_k\} = \{x_{\bar{k}+k}\}$. □

**Lemma 8.2.** *Let $x$ be a proper convex combination of atoms in $A' \subset A$, and $d \neq 0$ feasible direction in $x$. Then, for some $y \in \mathrm{conv}(A')$, we have*

$$
\hat{\alpha}^{\max}(y, d) \geq \frac{\mathrm{PWidth}(A)}{\|d\|}.
\tag{8.41}
$$

*Proof.* Let $y \in \mathrm{argmax}_{z \in \mathrm{conv}(A')} \hat{\alpha}^{\max}(z, d)$, and let $A'' \subset A'$ be such that $y$ is a proper convex combination of elements in $A''$. Furthermore, let $\mathcal{F}_y$ be the minimal face containing the maximal feasible step point $\bar{y} := y + \hat{\alpha}^{\max}(y, d)$. We claim that $\mathcal{F}_y \cap A'' = \emptyset$. In fact, for $p \in A'' \cap \mathcal{F}_y$ we can consider an homothety of center $p$ and factor $1 + \epsilon$ mapping $y$ in $y_\epsilon \in \mathrm{conv}(A'')$ and $\bar{y}$ in $\bar{y}_\epsilon \in \mathcal{F}_y$ with

$$
\bar{y}_\epsilon = y_\epsilon + (1 + \epsilon)\hat{\alpha}^{\max}(y, d)d.
$$

But then we would have $\hat{\alpha}(\bar{y}_\epsilon, d) \geq (1 + \epsilon)\hat{\alpha}(\bar{y}, d)$, in contradiction with the maximality of $\hat{\alpha}(\bar{y}, d)$. Therefore

$$
\hat{\alpha}^{\max}(y, d) \geq \mathrm{dist}(A'', \mathcal{F}_y) \geq \min_{\mathcal{F} \in \mathrm{pfaces}(\Omega)} \mathrm{dist}(\mathcal{F}, \mathrm{conv}(A \setminus \mathcal{F})) = \mathrm{PWidth}(A),
\tag{8.42}
$$

where we used $A'' \cap \mathcal{F} = \emptyset$ in the second inequality, and [55, Theorem 2] in the equality. □

31

# References

[1] P-A Absil, Robert Mahony, and Benjamin Andrews. Convergence of the iterates of descent methods for analytic cost functions. *SIAM Journal on Optimization*, 16(2):531–547, 2005.

[2] P-A Absil and Jérôme Malick. Projection-like retractions on matrix manifolds. *SIAM Journal on Optimization*, 22(1):135–158, 2012.

[3] Ralph Alexander. The width and diameter of a simplex. *Geometriae Dedicata*, 6(1):87–94, 1977.

[4] Hédy Attouch, Jérôme Bolte, Patrick Redont, and Antoine Soubeyran. Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the Kurdyka-Łojasiewicz inequality. *Mathematics of Operations Research*, 35(2):438–457, 2010.

[5] Hedy Attouch, Jérôme Bolte, and Benar Fux Svaiter. Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward–backward splitting, and regularized Gauss–Seidel methods. *Mathematical Programming*, 137(1-2):91–129, 2013.

[6] M. V. Balashov, B. T. Polyak, and A. A. Tremba. Gradient projection and conditional gradient methods for constrained nonconvex minimization. *Numerical Functional Analysis and Optimization*, 41(7):822–849, 2020.

[7] Mohammad Ali Bashiri and Xinhua Zhang. Decomposition-invariant conditional gradient for general polytopes with line search. In *Advances in Neural Information Processing Systems*, pages 2690–2700, 2017.

[8] Amir Beck and Shimrit Shtern. Linearly convergent away-step conditional gradient for non-strongly convex functions. *Mathematical Programming*, 164(1-2):1–27, 2017.

[9] Leonard Berrada, Andrew Zisserman, and M Pawan Kumar. Deep Frank-Wolfe for neural network optimization. In *International Conference on Learning Representations*, 2018.

[10] Dimitri P Bertsekas and Athena Scientific. *Convex optimization algorithms*. Athena Scientific Belmont, Nashua, 2015.

[11] Jérôme Bolte, Aris Daniilidis, Adrian Lewis, and Masahiro Shiota. Clarke subgradients of stratifiable functions. *SIAM Journal on Optimization*, 18(2):556–572, 2007.

[12] Jérôme Bolte, Aris Daniilidis, Olivier Ley, and Laurent Mazet. Characterizations of Łojasiewicz inequalities: subgradient flows, talweg, convexity. *Transactions of the American Mathematical Society*, 362(6):3319–3363, 2010.

[13] Jérôme Bolte, Trong Phong Nguyen, Juan Peypouquet, and Bruce W Suter. From error bounds to the complexity of first-order descent methods for convex functions. *Mathematical Programming*, 165(2):471–507, 2017.

[14] Jérôme Bolte, Shoham Sabach, and Marc Teboulle. Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Mathematical Programming*, 146(1):459–494, 2014.

[15] Immanuel M Bomze. Evolution towards the maximum clique. *Journal of Global Optimization*, 10(2):143–164, 1997.

[16] Immanuel M Bomze, Francesco Rinaldi, and Samuel Rota Bulo. First-order methods for the impatient: Support identification in finite time with convergent Frank-Wolfe variants. *SIAM Journal on Optimization*, 29(3):2211–2226, 2019.

[17] Immanuel M Bomze, Francesco Rinaldi, and Damiano Zeffiro. Active set complexity of the away-step Frank–Wolfe algorithm. *SIAM Journal on Optimization*, 30(3):2470–2500, 2020.

[18] Immanuel M Bomze, Francesco Rinaldi, and Damiano Zeffiro. Frank–wolfe and friends: a journey into projection-free first-order optimization methods. *4OR*, 19(3):313–345, 2021.

[19] Gábor Braun, Sebastian Pokutta, Dan Tu, and Stephen Wright. Blended conditonal gradients. In *International Conference on Machine Learning*, pages 735–743. PMLR, 2019.

[20] Gábor Braun, Sebastian Pokutta, and Daniel Zink. Lazifying conditional gradient algorithms. In *ICML*, pages 566–575, 2017.

[21] James V Burke and Jorge J Moré. On the identification of active constraints. *SIAM Journal on Numerical Analysis*, 25(5):1197–1211, 1988.

[22] Michael D Canon and Clifton D Cullum. A tight upper bound on the rate of convergence of Frank-Wolfe algorithm. *SIAM Journal on Control*, 6(4):509–516, 1968.

[23] Cyrille W Combettes and Sebastian Pokutta. Boosting Frank-Wolfe by chasing gradients. *arXiv preprint arXiv:2003.06369*, 2020.

[24] Andrea Cristofari, Marianna De Santis, Stefano Lucidi, and Francesco Rinaldi. An active-set algorithmic framework for non-convex optimization problems over the simplex. *Computational Optimization and Applications*, 77:57–89, 2020.

[25] Marguerite Frank and Philip Wolfe. An algorithm for quadratic programming. *Naval research logistics quarterly*, 3(1-2):95–110, 1956.

[26] Robert M Freund, Paul Grigas, and Rahul Mazumder. An extended Frank-Wolfe method with in-face directions, and its application to low-rank matrix completion. *SIAM Journal on Optimization*, 27(1):319–346, 2017.

[27] Dan Garber. Revisiting Frank-Wolfe for polytopes: Strict complementary and sparsity. *arXiv preprint arXiv:2006.00558*, 2020.

[28] Dan Garber and Ofer Meshi. Linear-memory and decomposition-invariant linearly convergent conditional gradient algorithm for structured polytopes. *Advances in neural information processing systems*, 29, 2016.

[29] Luigi Grippo, Francesco Lampariello, and Stephano Lucidi. A nonmonotone line search technique for newton's method. *SIAM Journal on Numerical Analysis*, 23(4):707–716, 1986.

[30] Peter Gritzmann and Marek Lassak. Estimates for the minimal width of polytopes inscribed in convex bodies. *Discrete & Computational Geometry*, 4(6):627–635, 1989.

[31] Jacques Guelat and Patrice Marcotte. Some comments on Wolfe's away step. *Mathematical Programming*, 35(1):110–119, 1986.

[32] Elad Hazan and Haipeng Luo. Variance-reduced and projection-free stochastic optimization. In *International Conference on Machine Learning*, pages 1263–1271, 2016.

[33] Martin Jaggi. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In *Proceedings of the 30th international conference on machine learning*, pages 427–435, 2013.

[34] Rujun Jiang and Xudong Li. Hölderian error bounds and kurdyka-łojasiewicz inequality for the trust region subproblem. *Mathematics of Operations Research*, 2022.

[35] Carl Johnell and Morteza Haghir Chehreghani. Frank-Wolfe optimization for dominant set clustering. *arXiv preprint arXiv:2007.11652*, 2020.

[36] David S Johnson. Cliques, coloring, and satisfiability: second dimacs implementation challenge. *DIMACS series in discrete mathematics and theoretical computer science*, 26:11–13, 1993.

[37] Armand Joulin, Kevin Tang, and Li Fei-Fei. Efficient image and video co-localization with Frank-Wolfe algorithm. In *European Conference on Computer Vision*, pages 253–268. Springer, 2014.

[38] Zoran Kadelburg, Dusan Dukic, Milivoje Lukic, and Ivan Matic. Inequalities of Karamata, Schur and Muirhead, and some applications. *The Teaching of Mathematics*, 8(1):31–45, 2005.

[39] Jovan Karamata. Sur une inégalité relative aux fonctions convexes. *Publications de l'Institut Mathématique*, 1(1):145–147, 1932.

[40] Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal-gradient methods under the Polyak-Łojasiewicz condition. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 795–811. Springer, 2016.

[41] Thomas Kerdreux, Alexandre d'Aspremont, and Sebastian Pokutta. Restarting Frank-Wolfe. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1275–1283. PMLR, 2019.

[42] Tamara G Kolda, Robert Michael Lewis, and Virginia Torczon. Stationarity results for generating set search for linearly constrained optimization. *SIAM Journal on Optimization*, 17(4):943–968, 2007.

[43] Vladimir Kolmogorov. Practical Frank-Wolfe algorithms. *arXiv preprint arXiv:2010.09567*, 2020.

[44] Simon Lacoste-Julien. Convergence rate of Frank-Wolfe for non-convex objectives. *arXiv preprint arXiv:1607.00345*, 2016.

[45] Simon Lacoste-Julien and Martin Jaggi. On the global linear convergence of frank-wolfe optimization variants. *Advances in neural information processing systems*, 28:496–504, 2015.

[46] Guanghui Lan and Yi Zhou. Conditional gradient sliding for convex optimization. *SIAM Journal on Optimization*, 26(2):1379–1409, 2016.

[47] Kfir Levy and Andreas Krause. Projection free online learning over smooth sets. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1458–1466, 2019.

[48] Robert Michael Lewis, Anne Shepherd, and Virginia Torczon. Implementing generating set search methods for linearly constrained minimization. *SIAM Journal on Scientific Computing*, 29(6):2507–2530, 2007.

[49] Guoyin Li and Ting Kei Pong. Calculus of the exponent of Kurdyka-Łojasiewicz inequality and its applications to linear convergence of first-order methods. *Foundations of computational mathematics*, 18(5):1199–1232, 2018.

[50] Stanislaw Lojasiewicz. Une propriété topologique des sous-ensembles analytiques réels. *Les équations aux dérivées partielles*, 117:87–89, 1963.

[51] Zhi-Quan Luo and Paul Tseng. Error bounds and convergence analysis of feasible descent methods: a general approach. *Annals of Operations Research*, 46(1):157–178, 1993.

[52] Hassan Mortagy, Swati Gupta, and Sebastian Pokutta. Walking in the shadow: A new perspective on descent directions for constrained minimization. *Advances in Neural Information Processing Systems*, 33, 2020.

[53] Julie Nutini, Mark Schmidt, and Warren Hare. "active-set complexity" of proximal gradient: How long does it take to find the sparsity pattern? *Optimization Letters*, 13(4):645–655, 2019.

[54] Anton Osokin, Jean-Baptiste Alayrac, Isabella Lukasewitz, Puneet Dokania, and Simon Lacoste-Julien. Minding the gaps for block frank-wolfe optimization of structured svms. In *International Conference on Machine Learning*, pages 593–602. PMLR, 2016.

[55] Javier Peña and Daniel Rodriguez. Polytope conditioning and linear convergence of the Frank-Wolfe algorithm. *Math. Oper. Res.*, 44(1):1–18, 2018.

[56] Fabian Pedregosa, Geoffrey Negiar, Armin Askari, and Martin Jaggi. Linearly convergent Frank-Wolfe with backtracking line-search. In *International Conference on Artificial Intelligence and Statistics*, pages 1–10. PMLR, 2020.

[57] Boris T Polyak. Gradient methods for the minimisation of functionals. *USSR Computational Mathematics and Mathematical Physics*, 3(4):864–878, 1963.

[58] Chao Qu, Yan Li, and Huan Xu. Non-convex conditional gradient sliding. In *International Conference on Machine Learning*, pages 4208–4217. PMLR, 2018.

[59] Luis Rademacher and Chang Shu. The smoothed complexity of Frank-Wolfe methods via conditioning of random matrices and polytopes. *arXiv preprint arXiv:2009.12685*, 2020.

[60] Francesco Rinaldi and Damiano Zeffiro. A unifying framework for the analysis of projection-free first-order methods under a sufficient slope condition. *arXiv preprint arXiv:2008.09781*, 2020.

[61] R Tyrrell Rockafellar and Roger J-B Wets. *Variational analysis*, volume 317. Springer Science & Business Media, Berlin, 2009.

[62] Arie Tamir. A strongly polynomial algorithm for minimum convex separable quadratic cost flow problems on two-terminal series-parallel networks. *Math. Program.*, 59:117–132, 1993.

[63] Klaus Truemper. Unimodular matrices of flow problems with additional constraints. *Networks*, 7(4):343–358, 1977.

[64] Yu Wang, Wotao Yin, and Jinshan Zeng. Global convergence of admm in non-convex nonsmooth optimization. *Journal of Scientific Computing*, 78(1):29–63, 2019.

[65] Philip Wolfe. Convergence theory in nonlinear programming. *Integer and nonlinear programming*, pages 1–36, 1970.

[66] Yangyang Xu and Wotao Yin. A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion. *SIAM Journal on imaging sciences*, 6(3):1758–1789, 2013.

[67] Li Zhang, Weijun Zhou, and Dong-Hui Li. A descent modified Polak-Ribière-Polyak conjugate gradient method and its global convergence. *IMA Journal of Numerical Analysis*, 26(4):629–640, 2006.