# RESTARTING FRANK-WOLFE: FASTER RATES UNDER HÖLDERIAN ERROR BOUNDS

THOMAS KERDREUX [†,∗], ALEXANDRE D'ASPREMONT[‡,§], AND SEBASTIAN POKUTTA[†,∗]

ABSTRACT. Conditional Gradient algorithms (aka Frank-Wolfe algorithms) form a classical set of methods for constrained smooth convex minimization due to their simplicity, the absence of projection steps, and competitive numerical performance. While the vanilla Frank-Wolfe algorithm only ensures a worst-case rate of $\mathcal{O}(1/\epsilon)$, various recent results have shown that for strongly convex functions on polytopes, the method can be slightly modified to achieve linear convergence. However, this still leaves a huge gap between sublinear $\mathcal{O}(1/\epsilon)$ convergence and linear $\mathcal{O}(\log 1/\epsilon)$ convergence to reach an $\epsilon$-approximate solution. Here, we present a new variant of Conditional Gradient algorithms, that can dynamically adapt to the function's geometric properties using restarts and smoothly interpolates between the sublinear and linear regimes. These interpolated convergence rates are obtained when the optimization problem satisfies a new type of error bounds, which we call *strong Wolfe primal bounds*. They combine geometric information on the constraint set with Hölderian Error Bounds on the objective function.

## 1. INTRODUCTION

We consider smooth constrained convex minimization, solving problems of the form

$$\min_{x \in \mathcal{C}} f(x), \tag{1}$$

where $f$ is a smooth convex function and $\mathcal{C}$ is a compact convex set. As soon as the geometry of $\mathcal{C}$ is reasonably complicated, so that projections onto the set are computationally expensive, projection-free first-order methods such as Conditional Gradient algorithms [34] (also known as Frank-Wolfe methods [14]) become an efficient alternative as they only require first-order access to the function under consideration as well as access to an efficient linear optimization oracle for the feasible region $\mathcal{C} \subseteq \mathbb{R}^n$ which, given a linear objective $c \in \mathbb{R}^n$, outputs $\arg\min_{x \in \mathcal{C}} c^T x$.

In order to reach an $\epsilon$-approximate solution $\hat{x}$, so that $f(\hat{x}) - f(x^*) < \epsilon$, where $x^*$ is an optimal solution, the standard Frank-Wolfe algorithm requires a number of iterations of order $O(1/\epsilon)$, that cannot be improved in general [9, 24]. A series of recent works (see e.g., [19, 30]; see also [32] for conditional gradient sliding) showed that when $f$ is strongly convex the convergence rate of the standard case can be improved to $O(\log 1/\epsilon)$ and various extensions further improved upon these results for special cases (see e.g., [31, 18, 16, 20, 7, 33, 3, 27, 8, 47, 12]), applying Frank-Wolfe methods to machine learning problems (e.g., [25, 46, 43, 36, 17, 37, 40]). Nonetheless, these results left a wide gap between the linear $O(\log 1/\epsilon)$ rate and the sub-linear $O(1/\epsilon)$ rate. Note that [29] also obtain such interpolated rates when analysing the vanilla Frank-Wolfe algorithm on (locally) uniformly convex constraint sets, extending the known accelerated regimes of Frank-Wolfe when the constraint set is strongly convex.

Here, we present a new variant of the Conditional Gradient method using the scaling argument of the parameter-free Lazy Frank-Wolfe variant in [7, 8], together with a restart scheme similar to that used for gradient methods in e.g., [41, 21, 42, 13, 45]. This yields an algorithm that dynamically adapts to the local properties of the function and the feasible region around the optimum. The convergence proof relies on two key conditions. One is a scaling inequality (Definition 3.3) used to characterize the regularity of $\mathcal{C}$ in many Frank-Wolfe complexity bounds which holds on e.g., polytopes and strongly convex sets. The other is a

local growth condition which is known to hold generically for sub-analytic functions by the Łojasiewicz factorization lemma (see e.g., [5]) and controls for example the impact of restart schemes as in [45].

Earlier work showed that a sharpness condition derived from the Łojasiewicz lemma could be used to improve convergence rates of gradient methods (see e.g., [41, 5, 1, 4, 26] for an overview). However, to achieve improved rates, these methods required exact knowledge of the constants appearing in the condition, which are in practice typically not observed. In contrast to this, as in [45, 10], we show using robust restart schemes that our algorithm does not require knowledge of these constants, thus making it essentially parameter-free.

**Contributions.** This paper is a journal version of [28]. Our contributions can be summarized as follows.

(1) *Strong Wolfe primal bound.* Under generic assumptions, we derive strong Wolfe primal gap bounds generalizing those obtained from strong convexity of $f$. These bounds are obtained by combining a Łojasiewicz growth condition on $f$ with a scaling inequality on $\mathcal{C}$, and continuously interpolate between the convex and strongly convex cases. Section 2 and 3 provide a more in-depth approach than in [28].

(2) *Fractional Frank-Wolfe Algorithms.* We then define a new Conditional Gradient algorithm that dynamically adapts to the parameters of these strong Wolfe primal bounds using a restart scheme. The resulting algorithm achieves either sub-linear (i.e., $O(1/\epsilon^q)$ with $q \leq 1$) or linear convergence rates depending on the strong Wolfe primal gap parameters. The exponent $q$ depends on the growth of the function around the optimum, so the function is not required to be strongly convex in the traditional sense. In particular, we obtain linear rates (depending on the parameters) for non-strongly convex functions. Our rates are satisfied after a mild burn-in phase that does not depend on the target accuracy. We extend the results of [28] to all the known settings where versions of Frank-Wolfe enjoy linear convergence rate under strong convexity assumptions of the objective function (see Section 7).

(3) *Robust restarts.* Restart schedules often heavily depend on the value of unknown parameters. We show that because Frank-Wolfe methods naturally produce a stopping criterion in the form of the strong Wolfe gap, our restart schemes are robust and do not require knowledge of the unobserved strong Wolfe primal gap bound parameters.

(4) We generalize our approach in [28] to Hölder smooth functions.

**Outline.** In Section 2 we briefly recall key notions and notation. We then describe our strong Wolfe primal bounds in Section 3 and present the Fractional Away-step Frank-Wolfe Algorithm in Section 4 along with the associated restart schemes in Section 5. Section 6 generalizes the analysis to Hölder smooth functions. Section 7 investigates the cases where the constraint set is not a polytope or the optimum is not necessarily on the boundary of the constraint set. They are known cases where additional structure on $f$ leads to accelerated convergence rates of the (vanilla) Frank-Wolfe algorithm.

## 2. PRELIMINARIES

Consider the following optimization problem

$$\begin{aligned}
\text{minimize} \quad & f(x) \\
\text{subject to} \quad & x \in \mathcal{C}
\end{aligned} \tag{2}$$

in the variables $x \in \mathbb{R}^n$, where $\mathcal{C} \subset \mathbb{R}^n$ is a compact convex set and $f : \mathbb{R}^n \to \mathbb{R}$ is a convex function. For the sake of simplicity, we will consider that $\mathcal{C}$ is full-dimensional. Let $X^*$ be the set of minimizers of $f$ over $\mathcal{C}$. We assume that the following linear minimization oracle

$$\text{LP}_{\mathcal{C}}(x) \triangleq \operatorname*{argmin}_{z \in \mathcal{C}} x^T z \tag{3}$$

can be computed efficiently. By assumption here, we have $\mathcal{C} = \mathbf{Co}(\mathbf{Ext}(\mathcal{C}))$ where $\mathbf{Co}(\cdot)$ is the convex hull, $\mathbf{Ext}(\cdot)$ the set of extreme points, and Carathéodory's theorem shows that every point $x$ of $\mathcal{C}$ can be

written as a convex combination of at most $n + 1$ points in $\mathbf{Ext}(\mathcal{C})$ although a given representation can contain more points. We call these points the *support of $x$* in $\mathcal{C}$. We say that a support $\mathcal{S}$ is *proper* when the weights that compose the convex combination of $x$ are all positive.

**Definition 2.1** (Proper Support). *Consider a compact convex set $\mathcal{C}$ and $x \in \mathcal{C}$. A finite set $\mathcal{S} = \{v_i \mid i \in I\}$ with $v_i \in \mathbf{Ext}(\mathcal{C})$ for some finite index set $I$, is a proper support of $x$ iff*

$$x = \sum_{i \in \mathcal{S}} \lambda_i v_i, \quad \text{where } \mathbf{1}^T \lambda = 1 \text{ and } \lambda_i > 0 \text{ for all } i \in I.$$

We now define the *strong Wolfe gap* as follows.

**Definition 2.2** (Strong Wolfe Gap). *Let $f$ be a smooth convex function, $\mathcal{C}$ a polytope, and let $x \in \mathcal{C}$ be arbitrary. Then the* strong Wolfe gap $w(x)$ *over $\mathcal{C}$ is defined as*

$$w(x) \triangleq \min_{\mathcal{S} \in \mathcal{S}_x} \max_{y \in \mathcal{S}, z \in \mathcal{C}} \nabla f(x)^T (y - z) \tag{4}$$

*where $x \in \mathbf{Co}(\mathcal{S})$ and $\mathcal{S}_x = \{\mathcal{S} \mid \mathcal{S} \subset \mathbf{Ext}(\mathcal{C}), \text{ is finite and } x \text{ a proper combination of the elements of } \mathcal{S}\}$, the set of proper supports of $x$. We also write*

$$w(x, \mathcal{S}) \triangleq \max_{y \in \mathcal{S}, z \in \mathcal{C}} \nabla f(x)^T (y - z)$$

*given $\mathcal{S} \in \mathcal{S}_x$.*

By construction, we have $w(x) \leq w(x, \mathcal{S})$. Note also that for $x \in \mathcal{C}$, the quantity $w(x, \mathcal{S})$ is the sum of the Frank-Wolfe dual gap with the away dual gap in [30] as shows the following decomposition

$$w(x, \mathcal{S}) = \underbrace{\max_{y \in \mathcal{S}} \nabla f(x)^T (y - x)}_{\text{away or Wolfe (dual) gap}} + \underbrace{\max_{z \in \mathcal{C}} \nabla f(x)^T (x - z)}_{\text{Frank-Wolfe (dual) gap}}. \tag{5}$$

Note that only $w(x, \mathcal{S})$ is observed in practice, but we use $w(x)$ to simplify the primal bounds and the convergence proof. Also we write the Frank-Wolfe (dual) gap as

$$g(x) \triangleq \max_{z \in \mathcal{C}} \nabla f(x)^T (x - z). \tag{6}$$

We first show the following lemma on $w(x, \mathcal{S})$ and $w(x)$. This lemma justifies the use of strong Wolfe gaps as measure of optimality in Algorithm 1 and 2.

**Lemma 2.3.** *Let $x \in \mathcal{C}$ and $\mathcal{S}$ be proper support of $x$. We have that $w(x, \mathcal{S}) = 0$ if and only if $x$ is an optimal solution of problem (2). In particular, $w(x) = 0$ if and only if $x$ is an optimal solution of problem (2). Write $x^*$ and optimal solution to (1), we have*

$$f(x) - f(x^*) \leq w(x, \mathcal{S}). \tag{7}$$

*Proof.* We can split $w(x, \mathcal{S})$ in two parts, with

$$w(x, \mathcal{S}) = \max_{y \in \mathcal{S}} \nabla f(x)^T (y - x) + \max_{z \in \mathcal{C}} \nabla f(x)^T (x - z).$$

If $x \in \mathcal{C}$, then both summands are nonnegative. Recall $g(x) = \max_{z \in \mathcal{C}} \nabla f(x)(x - z)$ is the Wolfe gap.

Let us first assume that $x$ is an optimal solution of problem (2) and show that $w(x, \mathcal{S}) = 0$. The first order optimality conditions implies that $\nabla f(x)^T (x - v) \leq 0$ for all $v \in \mathcal{C}$. Since this last quantity is exactly zero when $v = x$, we have $g(x) = 0$. Besides, let $h(x) \triangleq \max_{y \in \mathcal{S}} \nabla f(x)^T (y - x)$. If $\nabla f(x) = 0$ we immediately get $h(x) = 0$. Suppose then $\nabla f(x) \neq 0$, since $x$ is optimal, $\nabla f(x)^T (x - v_i) \leq 0$ for all $v_i \in \mathcal{S}$ and we can write

$$x = \sum_{\{i: \nabla f(x)^T (x - v_i) = 0\}} \lambda_i v_i + \sum_{\{i: \nabla f(x)^T (x - v_i) < 0\}} \lambda_i v_i = (1 - \mu) z_1 + \mu z_2$$

3

for some $0 \le \mu \le 1$, where $\nabla f(x)^T(x - z_1) = 0$ and $\nabla f(x)^T(x - z_2) < 0$. Now $0 = \nabla f(x)^T(x - x) = \mu \nabla f(x)^T(x - z_2)$ implies $\mu = 0$, hence $\nabla f(x)^T(x - v_i) = 0$ for all $i \in \mathcal{S}$, so $h(x) = 0$. Thus we obtain that $x$ optimal implies $w(x) = 0$.

Conversely, let us assume that $w(x, \mathcal{S}) = 0$. we have

$$
\begin{aligned}
f(x) - f^\star &\le \nabla f(x)^T(x - x^\star) \\
&\le \max_{z \in \mathcal{C}} \nabla f(x)^T(x - z) \\
&\le \max_{y \in S, z \in \mathcal{C}} \nabla f(x)^T(y - z) \\
&= w(x, \mathcal{S})
\end{aligned}
$$

by convexity (where $x^\star$ is any optimal solution), and the fact that $x \in \mathbf{Co}(\mathcal{S})$. Hence $w(x, \mathcal{S}) = 0$ implies $x$ optimal. The corollary on $w(x)$ immediately follows by construction. $\blacksquare$

A function $f$ is $L$-smooth when for any $x, y \in \mathcal{C}$

$$\|f(x) - f(y)\|_* \le L\|x - y\|. \tag{8}$$

Such regularity of the gradient can also be captured via curvature. We recall the definition of *away curvature* in [30, Appendix D], with

$$C_f^A \triangleq \sup_{\substack{x,s,v \in \mathcal{C} \\ \eta \in [0,1] \\ y = x + \eta(s-v)}} \frac{2}{\eta^2} \big(f(y) - f(x) - \eta\langle \nabla f(x), s - v \rangle\big), \tag{9}$$

where $f$ and $\mathcal{C}$ are defined in problem (2) above; we will use this notion of curvature for analyzing those algorithms utilizing away steps (Algorithm 1). Similarly (standard) *curvature* $C_f$ [30, Appendix C] is defined as

$$C_f \triangleq \sup_{\substack{x,v \in \mathcal{C} \\ \eta \in [0,1] \\ y = x + \eta(v-x)}} \frac{2}{\eta^2} \big(f(y) - f(x) - \eta\langle \nabla f(x), v - x \rangle\big), \tag{10}$$

and is used to bound the complexity of the classical Frank-Wolfe method (Algorithm 3).

## 3. HÖLDERIAN ERROR BOUNDS

We now introduce growth conditions used to bound the complexity of our variant of the Frank-Wolfe algorithm when solving the constrained optimization problem in (1). Let $\mathcal{C}$ be a general compact convex set with non-empty interior. The following condition will be at the core of our complexity analysis.

**Definition 3.1** (Strong Wolfe primal bound)**.** *Let $K$ be a compact neighborhood of $X^*$ in $\mathcal{C}$, where $X^*$ is the set of solutions of the constrained optimization problem (1). A function $f$ satisfies a $r$-strong Wolfe primal bound on $K$, if and only if there exists $r \ge 1$ and $\mu > 0$ such that for all $x \in K$*

$$f(x) - f^* \le \mu w(x)^r, \tag{11}$$

*and $f^*$ its optimal value.*

In the next section, provided $f$ is a smooth convex function, we will show, for instance, that $r = 2$ above guarantees linear convergence of our variant of Away Frank-Wolfe. This 2-strong Wolfe primal bound holds notably when $f$ is strongly convex over a polytope, which corresponds to the linear convergence bound in [30], hence the following observation.

**Observation 3.2** ($f$ strongly convex and $\mathcal{C}$ a polytope)**.** *The results in [30, Theorem 8 in Eq (28)] show that when $f$ is strongly convex and $\mathcal{C}$ is a polytope then there exists $\mu_f^A > 0$ such that for all $x \in \mathcal{C}$*

$$f(x) - f^* \le \frac{w(x)^2}{2\mu_f^A}.$$

*In other words, condition* (3.1) *holds with $r = 2$ in this case.*

The fact that $w(x) = 0$ if and only if $f(x) = f^*$ means that, in principle, the Łojasiewicz factorization lemma [5, §3.2.] could be used to show that condition (11) holds generically but with unobservable parameters. These parameters are inherently hard to infer because (11) combines the properties of $f$ and $\mathcal{C}$, not distinguishing between the contribution of the function from that of the structure of the constrained set. This was our initial approach in [28] but proving the subanalyticity of $w(\cdot)$ is however non-trivial.

Hence, although (11) has an appealing succinct form, our results will rely on the combination of a more classical Hölderian error bound (in Definition 3.5) defined on $f$, and a *scaling inequality* (defined below in Definition 3.3), essentially driven by the structure of the set $\mathcal{C}$. The combination of these two inequalities leads to a $r$-strong Wolfe primal bound. We first state the *scaling inequality* relative to the strong Wolfe gap $w(x)$ that we will use in the context of the the away step variant of the Frank-Wolfe algorithm.

**Definition 3.3** ($\delta$-scaling)**.** *A convex set $\mathcal{C}$ satisfies a* scaling inequality *if there exists $\delta(\mathcal{C}) > 0$ such that for all $x \in \mathcal{C} \setminus X^*$ and all differentiable convex function $f$,*

$$w(x) \geq \delta(\mathcal{C}) \max_{x^* \in X^*} \left\langle \nabla f(x), \frac{x - x^*}{\|x - x^*\|} \right\rangle. \qquad \text{(Scaling)}$$

Here again, the strong Wolfe gap $w(x)$ is the minimum over all proper supports of $x$ of the scalar product of the (negative) gradient with the pairwise direction formed by the difference of the Frank-Wolfe vertex and the away vertex. Hence the $\delta$-scaling inequality compares the worst pairwise FW direction with the normalization of the direction $x^* - x$. Notably this condition is known to hold when $\mathcal{C}$ is a polytope, with [30] showing the following result (see also [44] for a simpler variant).

**Lemma 3.4** ([30])**.** *A polytope satisfies the $\delta$-scaling inequality with $\delta(\mathcal{C}) = PWidth(\mathcal{C})$, where $PWidth(\mathcal{C})$ is the pyramidal width [30, (9)].*

We now recall the definition of the Hölderian error bound (aka sharpness) for a function $f$ on problem (1) [23, 38, 39, 5] (see e.g., [45] for more detailed references).

**Definition 3.5** (Hölderian error bound (HEB))**.** *Consider a convex function $f$ and $K$ a compact neighborhood of $X^*$ in $\mathcal{C}$. For optimization problem* (1)*, $f$ satisfies a $(\theta, c)$-HEB on $K$ if there exists $\theta \in [0, 1/2]$ and $c > 0$ such that for all $x \in K$*

$$\min_{x^* \in X^*} \|x - x^*\| \leq c(f(x) - f^*)^\theta. \qquad \text{(HEB)}$$

The Hölderian error bound (HEB) locally quantifies the behavior of $f$ around the constrained optimum of problem (2). A similar condition was used to show improved convergence rates for unconstrained optimization in e.g., [41, 2, 15, 26, 6, 45, 35]. Note that strong convexity implies $(\theta, c)$-HEB with $\theta = 1/2$ so (HEB) can be seen as a generalization of strong convexity. Here, $\theta$ will allow us to interpolate between sub-linear and linear convergence rates.

Finally, we show that when Problem (1) satisfies both $\delta$-Scaling and $(\theta, c)$-HEB, the $(1 - \theta)^{-1}$-strong Wolfe primal bound in (11) holds.

**Lemma 3.6.** *Assume $f$ is a differentiable convex function satisfying $(\theta, c)$-HEB on $K$, and that $\mathcal{C}$ satisfies $\delta$-Scaling inequality. Then for all $x \in K$*

$$f(x) - f^* \leq \left(\frac{c}{\delta}\right)^r w(x)^r,$$

*with $r = \frac{1}{1-\theta}$ and $f^*$ the objective value at constrained optima.*

*Proof.* Assume we have $(\theta, c)$-HEB on $K$. For $x \in K \setminus X^*$, by convexity, with $\tilde{x} \in \operatorname{argmin}_{x^* \in X^*} \|x - x^*\|$

$$f(x) - f^* \leq \frac{\langle \nabla f(x), x - \tilde{x} \rangle}{\|x - \tilde{x}\|} \|x - \tilde{x}\|.$$

5

Hence applying $(\theta, c)$-HEB leads to

$$
\begin{aligned}
f(x) - f^* &\leq c \frac{\langle \nabla f(x), x - \tilde{x} \rangle}{\|x - \tilde{x}\|} \left( f(x) - f^* \right)^\theta \\
&\leq c \max_{x^* \in X^*} \frac{\langle \nabla f(x), x - x^* \rangle}{\|x - x^*\|} \left( f(x) - f^* \right)^\theta,
\end{aligned}
\tag{12}
$$

from which we obtain

$$
f(x) - f^* \leq c^{\frac{1}{1-\theta}} \max_{x^* \in X^*} \left( \frac{\langle \nabla f(x), x - x^* \rangle}{\|x - x^*\|} \right)^{\frac{1}{1-\theta}}.
$$

Combining this with the $\delta$-scaling inequality, we have

$$
f(x) - f^* \leq \left( \frac{c}{\delta} \right)^{\frac{1}{1-\theta}} w(x)^{\frac{1}{1-\theta}},
$$

and the desired result. ∎

In the next section, varying values of $r \in [1, 2[$ in (11) allow to produce sub-linear complexity bounds of the form $\mathcal{O}(1/\epsilon^{1/(2-r)})$, continuously interpolating between the known sub-linear $\mathcal{O}(1/\epsilon)$ and a linear convergence rate obtained with $r = 2$. For simplicity of exposition, we will always pick $K = \mathcal{C}$ in what follows. We also write $\textbf{Int}(\cdot)$ for the interior of a set.

## 4. THE FRACTIONAL AWAY-STEP FRANK-WOLFE ALGORITHM

In this section, we focus on the case where $\mathcal{C}$ is a polytope and $f$ a smooth convex function. This means, in particular, that condition (Scaling) holds. We now state the Fractional Away-step Frank-Wolfe method as Algorithm 1, a variant of the Away-step Frank-Wolfe algorithm, tailored for restarting.

---

**Algorithm 1** Fractional Away-step Frank-Wolfe Algorithm

---

**Input:** A smooth convex function $f$ with curvature $C_f^A$. Starting point $x_0 = \sum_{v \in \mathcal{S}_0} \alpha_0^v v \in \mathcal{C}$ with support
$\quad \mathcal{S}_0 \subset \textbf{Ext}(\mathcal{C})$. LP oracle (3) and schedule parameter $\gamma > 0$.
1:  $t := 0$
2:  **while** $w(x_t, \mathcal{S}_t) > e^{-\gamma} w(x_0, \mathcal{S}_0)$ **do**
3:  $\quad v_t := \text{LP}_\mathcal{C}(\nabla f(x_t))$ and $d_t^{FW} := v_t - x_t$
4:  $\quad s_t := \text{LP}_{\mathcal{S}_t}(-\nabla f(x_t))$ with $\mathcal{S}_t$ current active set and $d_t^{Away} := x_t - s_t$
5:  $\quad$ **if** $-\nabla f(x_t)^T d_t^{FW} > e^{-\gamma} w(x_0, \mathcal{S}_0)/2$ **then**
6:  $\quad\quad d_t := d_t^{FW}$ with $\eta_{\max} := 1$
7:  $\quad$ **else**
8:  $\quad\quad d_t := d_t^{Away}$ with $\eta_{\max} := \frac{\alpha_t^{s_t}}{1 - \alpha_t^{s_t}}$
9:  $\quad$ **end if**
10:  $\quad x_{t+1} := x_t + \eta_t d_t$ with $\eta_t \in [0, \eta_{\max}]$ via line-search
11:  $\quad$ Update active set $\mathcal{S}_{t+1}$ and coefficients $\{\alpha_{t+1}^v\}_{v \in \mathcal{S}_{t+1}}$
12:  $\quad t := t + 1$
13:  **end while**
**Output:** $t \in \mathbb{N}$ and $x_t \in \mathcal{C}$ such that $w(x_t, \mathcal{S}_t) \leq e^{-\gamma} w(x_0, \mathcal{S}_0)$

---

For the Away-step Frank-Wolfe or Algorithm 1, an iteration performs a *drop step* when the update direction is the away direction (Line 8) and the chosen step size $\eta_t$ is equal to $\eta_{\max} = \alpha_{s_t}/(1 - \alpha_{s_t})$. Indeed, such an iteration removes (drops) the away vertex $s_t$ from the convex combination of $x_t$. Conversely, we will call a step a *full-progress step* if it is a Frank-Wolfe step or an away step that is not a drop step. The support $\mathcal{S}_t$ and the weights $\alpha_t$ are updated exactly as in [30, Away-step Frank-Wolfe]. Note that to perform

6

such away versions of Frank-Wolfe, we require a linear minimization oracle over the support of the optimization iterates. Such an oracle is typically done naively so that its cost grows linearly with the size of the support. Algorithm 1 depends on a parameter $\gamma > 0$ which explicitly controls the number of iterations needed for the algorithm to stop. In particular, a large value of $\gamma$ will increase the number of iterations and when $\gamma$ converges to infinity, Algorithm 1 tends to behave exactly like the classical Frank-Wolfe and never chooses the away direction as an update direction. To support this intuition, we prove Appendix A that the convergence rate of one run of Fractional Away Frank-Wolfe with a large value of $\gamma$ similar to that of the classical Frank-Wolfe.

We name Algorithm 1 a *Fractional* version of Away Frank-Wolfe since after running the algorithm, the strong Wolfe gap $w(x_t, \mathcal{S}_t)$ (our measure of optimality) is only guaranteed to be a fraction of the initial Wolfe gap $w(x_0, \mathcal{S}_0)$. Besides, the vanilla AFW consists in a different decision rule to decide between away-steps or FW steps (Line 5).

Proposition 4.1 below gives an upper bound on the number of iterations required for Algorithm 1 to reach a given target gap $w(x_T, \mathcal{S}_T) \leq w(x_0, \mathcal{S}_0)e^{-\gamma}$. The assumption $e^{-\gamma}w(x_0, \mathcal{S}_0)/2 \leq C_f^A$ in this proposition measures the complexity of a burn-in phase whose cost is marginal as shown in Proposition 4.2.

**Proposition 4.1** (Fractional Away-Step Frank-Wolfe Complexity). *Let $f$ be a smooth convex function with away curvature $C_f^A$ such that the $r$-strong Wolfe primal bound in (11) holds on $\mathcal{C}$ (with $1 \leq r \leq 2$ and $\mu > 0$). Let $\gamma > 0$ and assume $x_0 \in \mathcal{C}$ is such that $e^{-\gamma}w(x_0)/2 \leq C_f^A$. Algorithm 1 outputs an iterate $x_T \in \mathcal{C}$ such that*

$$w(x_T, \mathcal{S}_T) \leq w(x_0, \mathcal{S}_0)e^{-\gamma}$$

*after at most*

$$T \leq |\mathcal{S}_0| - |\mathcal{S}_T| + 16e^{2\gamma}C_f^A \mu w(x_0, \mathcal{S}_0)^{r-2}$$

*iterations, where $\mathcal{S}_0$ and $\mathcal{S}_T$ are respectively the supports of $x_0$ and $x_T$.*

*Proof.* Let us write $w_0 \triangleq w(x_0, \mathcal{S}_0)$ to simplify notation. Because of the test criterion in Line 5, one can lower bound the inner product between the update directions $d_t$ and the negative gradients $-\nabla f(x_t)$ of the form

$$-\nabla f(x_t)^T d_t > e^{-\gamma}w(x_0, \mathcal{S}_0)/2 . \tag{13}$$

Indeed, this holds by definition when $d_t = d_t^{FW}$. Otherwise, $d_t = d_t^{Away}$ and $-\nabla f(x_t)^T d_t^{FW} < e^{-\gamma}w_0/2$. Also because the algorithm has not terminated yet, we have $w(x_t, \mathcal{S}_t) > e^{-\gamma}w_0$. The decomposition of the strong Wolfe gap (5) then yields

$$w(x_t, \mathcal{S}_t) = -\nabla f(x_t)^T d_t^{FW} - \nabla f(x_t)^T d_t^{Away} > e^{-\gamma}w_0,$$

so that we indeed obtain (13)

$$-\nabla f(x_t)^T d_t^{Away} > e^{-\gamma}w_0 + \nabla f(x_t)^T d_t^{FW} \geq e^{-\gamma}w_0 - e^{-\gamma}w_0/2 = e^{-\gamma}w_0/2.$$

Using curvature in (9), we have for $d_t$,

$$f(x_t + \eta d_t) \leq f(x_t) + \eta \nabla f(x_t)^T d_t + \frac{\eta^2}{2}C_f^A ,$$

which implies

$$f(x_t) - f(x_t + \eta d_t) \geq \eta - \nabla f(x_t)^T d_t - \frac{\eta^2}{2}C_f^A .$$

We can lower bound progress $f(x_t) - f(x_{t+1})$ with $x_{t+1} = x_t + \eta d_t$ at each iteration for full-progress steps. Indeed, for Frank-Wolfe steps,

$$
\begin{aligned}
f(x_t) - f(x_{t+1}) &\geq \max_{\eta \in [0,1]}\left\{\eta(-\nabla f(x_t))^T d_t - \frac{\eta^2}{2}C_f^A\right\} \\
&\geq \max_{\eta \in [0,1]}\left\{\eta e^{-\gamma}w_0/2 - \frac{\eta^2}{2}C_f^A\right\}
\end{aligned}
$$

7

Hence because of exact line-search, assuming $e^{-\gamma}w_0/2 \leq C_f^A$ holds, we obtain

$$f(x_t) - f(x_{t+1}) \geq \frac{w_0^2}{8C_f^A e^{2\gamma}}. \tag{14}$$

For all away steps, we have

$$f(x_t) - f(x_t + \eta d_t) \geq \max_{\eta \in [0, \eta_{\max}]} \left\{ \eta e^{-\gamma}w_0/2 - \frac{\eta^2}{2}C_f^A \right\}.$$

Yet for Away steps that are not drop steps, assuming $e^{-\gamma}w_0/2 \leq C_f^A$ holds, the minimal $\eta^*$ is such that $0 < \eta^* < \eta_{\max}$, and the same conclusion as in (14) for Frank-Wolfe steps follows.

Write $T = T_d + T_f$ the number of iterations for Algorithm 1 to finish, where $T_d$ denotes the number of drop steps, while $T_f$ stands for the number of full-progress steps. Hence we have,

$$
\begin{aligned}
f(x_0) - f(x_T) &= \sum_{t=0}^{T-1} f(x_t) - f(x_{t+1}) \\
&\geq T_f \frac{w_0^2}{8C_f^A e^{2\gamma}}.
\end{aligned}
$$

Because $f$ satisfies a $r$-strong Wolfe primal gap on $\mathcal{C}$ we have, when $x_0 \in \mathcal{C}$,

$$f(x_0) - f(x_T) \leq f(x_0) - f^* \leq \mu w(x_0)^r \leq \mu w(x_0, \mathcal{S}_0)^r, \tag{15}$$

by definition of $w(x)$. We then get an upper bound on the number $T_f$ of full-progress steps

$$T_f \leq 8C_f^A e^{2\gamma} \mu w_0^{r-2}.$$

Finally writing $|\mathcal{S}_0|$ (resp. $|\mathcal{S}_T|$) the size of the support of $x_0$ (resp. $x_T$), and $T_{FW}$ the number of Frank-Wolfe steps which add a new vertex to an iterate of the Fractional Away-step Frank-Wolfe Algorithm. We have that $T_{FW} \leq T_f$ and the size of the support $\mathcal{S}_t$ of $x_t$ satisfies $|\mathcal{S}_0| - T_d + T_{FW} = |\mathcal{S}_T|$ hence

$$|\mathcal{S}_0| - |\mathcal{S}_T| + T_f \geq T_d,$$

and we finally obtain $T \leq |\mathcal{S}_0| - |\mathcal{S}_T| + 16C_f^A e^{2\gamma} \mu w_0^{r-2}$. ∎

The following observation shows that the assumption $e^{-\gamma}w(x_0, \mathcal{S}_0)/2 \leq C_f^A$ in Proposition 4.1 has a marginal impact on complexity.

**Proposition 4.2** (Burn-in phase). *After at most*

$$8e^\gamma \left\lfloor \frac{1}{\gamma} \ln \frac{w_0}{2C_f^A} \right\rfloor + |\mathcal{S}_0|,$$

*cumulative iterations of Algorithm 1, with constant schedule parameter $\gamma > 0$, we obtain a point $x \in \mathcal{C}$ such that $e^{-\gamma}w(x, \mathcal{S}_x)/2 \leq C_f^A$.*

*Proof.* The proof closely follows that of Proposition 4.1 as well as [30, Appendix D]. Let $w_0 = w(x_0, \mathcal{S}_0)$. Suppose that $e^{-\gamma}w_0/2 > C_f^A$ and note that the lower bound (13) holds similarly. Let us consider the progress incurred with *full progress steps*. Recall that the curvature of $f$ with the line-search ensure that for any $\eta \in [0, \eta_{\max}]$, we have

$$f(x_t) - f(x_{t+1}) \geq \eta(-\nabla f(x_t))^T d_t - \frac{\eta^2}{2}C_f^A.$$

For a Frank-Wolfe step or an away step with $\eta_{\max} \geq 1$, we obtain by choosing $\eta = 1$,

$$f(x_t) - f(x_{t+1}) \geq (-\nabla f(x_t))^T d_t - C_f^A/2 \geq e^{-\gamma}w_0/2 - e^{-\gamma}w_0/4 = e^{-\gamma}w_0/4.$$

The last possibility is that the full progress step is an away step with $\eta_{\max} < 1$. Since it is not a drop step, we have $\eta_t < \eta_{\max}$. In particular then $\eta_t$ is a local minimum of the convex function $f(x_t + \gamma d_t)$ and hence $\min_{\gamma \in [0, \gamma_{\max}]} f(x_t + \gamma d_t) = \min_{\gamma} f(x_t + \gamma d_t)$. With $\eta = 1$, we obtain

$$f(x_t) - f(x_{t+1}) \geq (-\nabla f(x_t))^T d_t - C_f^A / 2.$$

Finally, we conclude that for any full progress step we have

$$f(x_t) - f(x_{t+1}) \geq e^{-\gamma} w_0 / 4.$$

Moreover, the strong Wolfe gap is an upper bound on the primal gap, *i.e.* $f(x_0) - f(x^*) \leq w_0$. Write $T$ the number of iterations Algorithm 1 performs. As in Proposition 4.1, note $T_f$ the number of full progress steps. Similarly, we obtain

$$T_f e^{-\gamma} w_0 / 4 \leq \sum_{t=0}^{T-1} f(x_t) - f(x_{t+1}) = f(x_0) - f(x_T) \leq f(x_0) - f(x^*) \leq w_0.$$

Hence

$$T_f e^{-\gamma} w_0 / 4 \leq w_0$$

and $T_f \leq 4e^{\gamma}$. Also

$$T = T_d + T_f \leq 2T_f + |\mathcal{S}_0| - |\mathcal{S}_T|,$$

so that

$$T \leq 8e^{\gamma} + |\mathcal{S}_0| - |\mathcal{S}_T|.$$

In other words, when $e^{-\gamma} w(x_0) / 2 < C_f^A$, after at most $8e^{\gamma} + |\mathcal{S}_0|$ iterations, the Fractional Away-step Frank-Wolfe terminates, with an iterate $x_T$ which strong Wolfe gap is guaranteed to be a fraction of the initial one, *i.e.*, $w(x_T, \mathcal{S}_T) \leq e^{-\gamma} w(x_0, \mathcal{S}_0)$.

Then, consider running the Fractional Away-step Frank-Wolfe $N$ times, initializing each run with the output of the previous run. Write $x_{T_i}$ each output of the $i^{th}$ run of Algorithm 1. After N runs, $x_{T_N}$ satisfies $w(x_{T_N}, \mathcal{S}_{T_N}) \leq e^{-\gamma N} w_0$. Hence, if $N$ satisfies

$$e^{-\gamma(N+1)} w_0 / 2 \leq C_f^A,$$

then $x_{T_N}$ verifies $e^{-\gamma} w(x_{T_N}, \mathcal{S}_{T_N}) \leq C_f^A$. In particular, it is sufficient to chose $N = \lfloor \frac{1}{\gamma} \ln w_0 / (2C_f^A) \rfloor$.

Finally, since the $i^{th}$ run of Fractional Away-step Frank-Wolfe performs at most $8e^{\gamma} + |\mathcal{S}_{T_{i-1}}| - |\mathcal{S}_{T_i}|$ iterations, to ensure the burn-in phase condition, we need at most

$$\sum_{i=1}^{N} 8e^{\gamma} + |\mathcal{S}_{T_{i-1}}| - |\mathcal{S}_{T_i}| \leq 8e^{\gamma} \left\lfloor \frac{1}{\gamma} \ln \frac{w_0}{2C_f^A} \right\rfloor + |\mathcal{S}_0| \text{ iterations.}$$

∎

## 5. RESTART SCHEMES

Consider a point $x_{k-1}$ with strong Wolfe gap $w(x_{k-1}, \mathcal{S}_{k-1})$. Algorithm 1 with parameter $\gamma_k > 0$, outputs a point $x_k$ and we write

$$x_k \triangleq \mathcal{F}(x_{k-1}, w(x_{k-1}, \mathcal{S}_{k-1}), \gamma_k).$$

Following [45] we define *scheduled restarts* for Algorithm 1 as follows.

**Algorithm 2** Scheduled restarts for Fractional Away-step Frank-Wolfe

---

**Input:** $\tilde{x}_0 \in \mathbb{R}^n$ and a sequence $(\gamma_k) > 0$ and $\epsilon > 0$ and $T \in \mathbb{N}$.
$t := 0$
**while** $w(x_{k-1}, \mathcal{S}_{k-1}) > \epsilon$ and $t < T$ **do**

$$(x_k, \mathcal{S}_k, T_k) := \mathcal{F}(x_{k-1}, w(x_{k-1}, \mathcal{S}_{k-1}), \gamma_k) \text{ and } t := t + T_k.$$

**end while**
**Output:** $x_k$

---

Note that one overall burn-in phase is sufficient to ensure the condition $e^{-\gamma_i} w(x_{i-1}, \mathcal{S}_{i-1})/2 \leq C_f^A$ at each restart.

Algorithm 2 is similar to the restart scheme in [45, Section 4] where a termination criterion is available. In this situation, [45] show that the convergence rate of restarted gradient methods is robust to a suboptimal choice of restart scheme parameter $\gamma$. Here, we also show that our restart scheme is adaptive to the unknown parameters in $(\theta, c)$-HEB.

Note that Algorithm 2 shares a similar structure with the methods in [33, 8]. We will see below in Proposition 5.3 that tuning $\gamma$ only has a marginal impact on the complexity bound. Note also that when $\theta \in [0, 1/2]$, the condition interpolates between the non-strongly convex function $f$ and a strongly convex function scenarios. For the sake of clarity, our convergence results depend on a burn-in phase condition on the initial strong Wolfe gap, *i.e.* $e^{-\gamma} w_0/2 \leq C_f^A$. Proposition 4.2 shows that it is satisfied after an initial linear convergence regime.

**Theorem 5.1** (Rate for constant restart schemes). *Let $f$ be a smooth convex function with away curvature $C_f^A$. Assume $\mathcal{C}$ satisfies $\delta$-Scaling and $f$ is $(\theta, c)$-HEB on $\mathcal{C}$. Let $\gamma > 0$ and assume $x_0 \in \mathcal{C}$ is such that $e^{-\gamma} w(x_0, \mathcal{S}_0)/2 \leq C_f^A$ (see, Proposition 4.2). With $\gamma_k = \gamma$, the output of Algorithm 2 satisfies ($r = \frac{1}{1-\theta}$)*

$$
\begin{cases}
f(x_T) - f^* \leq w_0 \dfrac{1}{\left(1 + \tilde{T} C_\gamma^r\right)^{\frac{1}{2-r}}} & \text{when } 1 \leq r < 2 \\[4mm]
f(x_T) - f^* \leq w_0 \exp\left(-\dfrac{\gamma}{e^{2\gamma}} \dfrac{\tilde{T}}{16 C_f^A \mu}\right) & \text{when } r = 2,
\end{cases}
\tag{16}
$$

*where $T$ is the cumulative number of Linear Minimization Oracle calls (Line 3 in Algorithm 1) in Algorithm 2, with $w_0 = w(x_0, \mathcal{S}_0)$, $\tilde{T} \triangleq T - (|\mathcal{S}_0| - |\mathcal{S}_T|)$, and*

$$
C_\gamma^r \triangleq \frac{e^{\gamma(2-r)} - 1}{16 C_f^A \mu e^{2\gamma} w(x_0, \mathcal{S}_0)^{r-2}},
\tag{17}
$$

*with $\mu = \frac{c}{\delta}$.*

*Proof.* Denote by $R$ the number of restarts in Algorithm 1 for $T$ total iterations. By design

$$
w(x_R, \mathcal{S}_R) \leq w_0 e^{-\gamma R}.
\tag{18}
$$

Because $f$ is $(\theta, c)$-HEB and $\mathcal{C}$ satisfies $\delta$-Scaling, via Lemma 3.6, $f$ satisfies the $r$-strong Wolfe primal bound (11) with $r = \frac{1}{1-\theta}$. Using Proposition 4.1 and repeatedly (18), the total number $T$ of steps of Algorithms 1 is upper-bounded by

$$
T \leq |\mathcal{S}_0| - |\mathcal{S}_T| + 16 C_f^A \mu e^{2\gamma} w_0^{r-2} \sum_{i=0}^{R-1} e^{-\gamma i (r-2)}.
$$

Let us now distinct the case and first suppose that $1 \leq r < 2$. We have the following upper bound on $T$,

$$
T \leq |\mathcal{S}_0| - |\mathcal{S}_T| + 16 C_f^A \mu e^{2\gamma} w_0^{r-2} \frac{e^{\gamma(2-r)R} - 1}{e^{\gamma(2-r)} - 1} = |\mathcal{S}_0| - |\mathcal{S}_T| + \frac{e^{\gamma(2-r)R} - 1}{C_\gamma^r},
$$

hence, with $\tilde{T} = T - |\mathcal{S}_0| + |\mathcal{S}_T|$,

$$e^{-\gamma R} \leq \frac{1}{\left(1 + \tilde{T}C_\gamma^r\right)^{\frac{1}{2-r}}}.$$

Thus, for $1 \leq r < 2$,

$$w(x_R, \mathcal{S}_R) \leq w_0 \frac{1}{\left(1 + \tilde{T}C_\gamma^r\right)^{\frac{1}{2-r}}}.$$

Now, the remaining case $r = 2$ leads to

$$T \leq |\mathcal{S}_0| - |\mathcal{S}_T| + 16C_f^A \mu e^{2\gamma} R,$$

and hence

$$w(x_R, \mathcal{S}_R) \leq w_0 \exp\left(-\gamma \frac{\tilde{T}}{16C_f^A \mu e^{2\gamma}}\right),$$

which yields the desired result. ∎

**Corollary 5.2.** *When $\mathcal{C}$ is a polytope and $f$ a smooth convex function satisfying $(\theta, c)$-HEB, rates in Theorem 5.1 hold. In particular when $f$ is strongly convex, $\theta = \frac{1}{2}$ (and hence $r = 2$) and Algorithm 2 converges linearly. When $f$ is simply smooth, $\theta = 0$ (and hence $r = 1$) and Algorithm 2 converges sub-linearly with a rate of $\mathcal{O}(1/T)$.*

Note also that for $r \to 2$, we recover the same complexity rates as for $r = 2$

$$\lim_{r \to 2} \frac{1}{\left(1 + \tilde{T}C_\gamma^r\right)^{\frac{1}{2-r}}} = \exp\left(-\frac{\gamma}{e^{2\gamma}} \frac{\tilde{T}}{8C_f^A \mu}\right).$$

The complexity bounds in Theorem 5.1 depend on $\gamma$, which controls the convergence rate. Optimal choices of $\gamma$ depend on $r$, a constant that we generally do not know nor observe. However, in the following we show that simply picking $\gamma = 1/2$ leads to optimal complexity bounds up to a constant factor. In fact, picking a constant gamma (independent of $r$) we also recover a simple version of [8, Algorithm 1] (without the cheaper Weak Separation Oracle that replaces the Linear Minimization Oracle).

**Proposition 5.3** (Robustness in $\gamma$). *Suppose $f$ satisfies the $r$-strong Wolfe primal bound (11) with $r > 0$. Write $\gamma^*(r)$ as the optimal choice of $\gamma > 0$ in the coarser complexity bounds (16) of Theorem 5.1 where $\tilde{T}$ is lower bounded by $\bar{T} = T - |\mathcal{S}_0|$. Consider running Algorithm 2 with $\gamma = 1/2$ and the same assumptions as in Theorem 5.1, the output $\hat{x}$ satisfies*

$$h(\hat{x}) \leq \sqrt{\frac{e}{4(\sqrt{e}-1)}} w_0 \frac{1}{\left(1 + \bar{T}C_{\gamma^*(r)}^r\right)^{\frac{1}{2-r}}} \quad \text{when } 1 \leq r < 2,$$

*where*

$$C_\gamma^r = \frac{e^{\gamma(2-r)} - 1}{16e^{2\gamma}C_f^A \mu w(x_0, \mathcal{S}_0)^{r-2}},$$

*as in (17). When $r = 2$, we have $\gamma^*(r) = 1/2$.*

*Proof.* When $1 \leq r < 2$, from Theorem 5.1 we have

$$f(x_T) - f^* \leq w_0 \frac{1}{\left(1 + \bar{T}C_\gamma^r\right)^{\frac{1}{2-r}}}. \tag{19}$$

With definition of $C_\gamma^r$ in (17), minimizing (19) is equivalent to maximizing (for $\gamma > 0$)

$$B(\gamma) = \left(\frac{e^{\gamma(2-r)} - 1}{e^{2\gamma}}\right).$$

11

Hence the optimum schedule parameter $\gamma^*(r)$ is

$$\gamma^*(r) = \frac{\ln(2) - \ln(r)}{2 - r} \quad \text{when } 1 \leq r < 2.$$

In particular $\gamma^*(r) \in ]1/2; \ln(2)]$. Let's now show that the bound in (19) obtained with the optimal $\gamma^*(r)$ is comparable to the bound obtained with $\gamma = \frac{1}{2}$. The function

$$H(r) = \frac{\left(1 + \bar{T}C^r_{\gamma^*(r)}\right)^{\frac{1}{2-r}}}{\left(1 + \bar{T}C^r_{1/2}\right)^{\frac{1}{2-r}}}$$

is decreasing in $r$. Write $\tilde{C} \triangleq 8C_f^A \mu w(x_0, \mathcal{S}_0)$, we have $C^1_{\gamma^*(1)} = 1/(4\tilde{C})$ and $C^1_{1/2} = \frac{\sqrt{e}-1}{e}/\tilde{C}$ hence

$$H(1) = \sqrt{\frac{1 + \frac{\bar{T}}{\tilde{C}}\frac{1}{4}}{1 + \frac{\bar{T}}{\tilde{C}}\frac{\sqrt{e}-1}{e}}} \leq \sqrt{\frac{e}{4(\sqrt{e}-1)}} \ .$$

Hence, with $H(1) \geq H(r)$, we get for any $r \in [1, 2[$

$$\frac{1}{\left(1 + \bar{T}C^r_{1/2}\right)^{\frac{1}{2-r}}} \leq \sqrt{\frac{e}{4(\sqrt{e}-1)}} \frac{1}{\left(1 + \bar{T}C^r_{\gamma^*(r)}\right)^{\frac{1}{2-r}}}.$$

When $r = 2$, the optimal choice for $\gamma$ is $1/2$, maximizing the function $\gamma/e^{2\gamma}$. ∎

In Figure 1, we illustrate the convergence behavior of Algorithm 2 along with that of the Away Frank-Wolfe. The algorithms have a similar behavior in the primal gap $f(x_t) - f(x^*)$. For numerically competitive *corrective* versions of Frank-Wolfe, see, e.g., [20, 3, 11] and references therein.
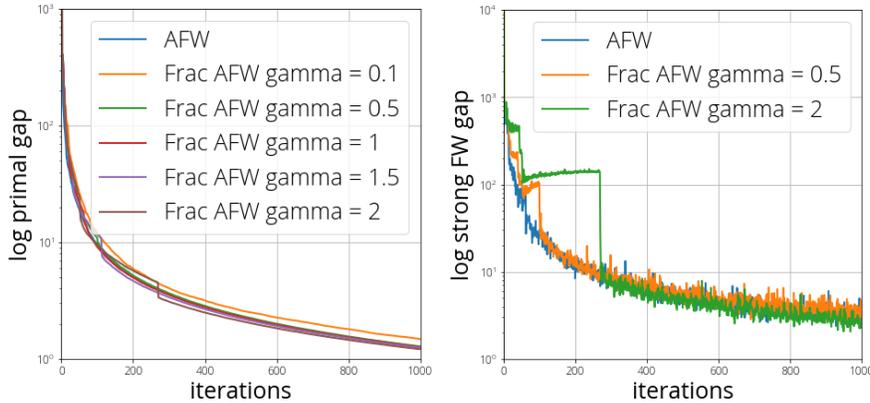


FIGURE 1. Representative examples on Lasso with various values of $\gamma$ in restart schemes of Algorithm 1.

## 6. ANALYSIS UNDER HÖLDER SMOOTHNESS

In the following we generalize our results on convergence rates using a refined regularity assumption on $f$. A differentiable function $f$ is $(L, s)$-Hölder smooth on $\mathcal{C}$ when

$$\|\nabla f(x) - \nabla f(y)\|_2 \leq L\|x - y\|_2^{s-1}, \quad \text{for } x, y \in \mathcal{C},$$

with $s \in ]1, 2]$. Hölder smoothness interpolates between non-smooth ($s = 1$) and smooth ($s = 2$) assumptions. We write the analog of the away curvature (9) for $(L, s)$-Hölder smooth functions as

$$C_{f,s}^A \triangleq \sup_{\substack{x,u,v \in \mathcal{C} \\ \eta \in [0,1] \\ y=x+\eta(u-v)}} \frac{s}{\eta^s} \big(f(y) - f(x) - \eta \langle \nabla f(x), u - v \rangle \big).$$

Note that as in (9), $f$ needs to be defined on the Minkowski sum $\mathcal{C}^A$. Let us now provide equivalent results for the complexity of Fractional Away-step Frank-Wolfe algorithm and the complexity bound of the constant restart scheme with $(L, s)$-Hölder smooth functions.

**Proposition 6.1** (Hölder Smooth Complexity). *Let $f$ be a $(L, s)$-Hölder smooth convex function with away curvature $C_{f,s}^A$ such that the $r$-strong Wolfe primal bound in (11) holds on $\mathcal{C}$ with $\mu > 0$. Let $\gamma > 0$ and assume $x_0 \in \mathcal{C}$ is such that $e^{-\gamma}w(x_0, \mathcal{S}_0)/2 \leq C_{f,s}^A$. Algorithm 1 outputs an iterate $x_T \in \mathcal{C}$ such that*

$$w(x_T, \mathcal{S}_T) \leq w(x_0, \mathcal{S}_0)e^{-\gamma}$$

*after at most (with $r = \frac{1}{1-\theta}$)*

$$T \leq |\mathcal{S}_0| - |\mathcal{S}_T| + 2^{1+\frac{s}{s-1}} \frac{s}{s-1} e^{\frac{s}{s-1}\gamma} \mu \big(C_{f,s}^A\big)^{\frac{1}{s-1}} w(x_0, \mathcal{S}_0)^{r-\frac{s}{s-1}}$$

*iterations, where $\mathcal{S}_0$ and $\mathcal{S}_T$ are the supports of respectively $x_0$ and $x_T$.*

*Proof.* The proof is very similar to that required for smooth-functions, so we only detail key points. The update direction satisfies

$$(-\nabla f(x_t))^T d_t > e^{-\gamma}w_0/2 .$$

Applying the definition of the Hölder curvature

$$f(x_t) - f(x_t + \eta d_t) \geq \max_{\eta \in [0,\eta_{\max}]} \{\eta e^{-\gamma}w_0/2 - \frac{\eta^s}{s}C_{f,s}^A\} = \max_{\eta \in [0,\eta_{\max}]} g(\eta) .$$

The unconstrained maximum of $g$ is reached at $\eta^* = \left(e^{-\gamma}\frac{w_0}{2C_{f,s}^A}\right)^{\frac{1}{s-1}}$. Hence with the burn-in phase hypothesis, we guarantee $\eta^* \leq 1$. With classical arguments, for all non-drop steps, the progress in the objective function value is lower bounded by

$$f(x_t) - f(x_t + \eta d_t) \geq \frac{1}{\left(C_{f,s}^A\right)^{\frac{1}{s-1}}} \frac{s-1}{s} 2^{-\frac{s}{s-1}} e^{-\gamma\frac{s}{s-1}} w_0^{\frac{s}{s-1}} .$$

It finally follows that

$$T \leq 2\mu w_0^{r-\frac{s}{s-1}} 2^{\frac{s}{s-1}} \frac{s}{s-1} e^{\gamma\frac{s}{s-1}} + |\mathcal{S}_0| - |\mathcal{S}_T|$$

which is the desired bound. $\blacksquare$

We are ready to establish the convergence rates of our restart scheme in the Hölder smooth case.

**Theorem 6.2** (Hölder rate for constant restart schemes). *Let $f$ be a $(L, s)$-Hölder smooth convex function with Hölder curvature $C_{f,s}^A$, satisfying $(\theta, c)$-HEB on $\mathcal{C}$, and $\mathcal{C}$ satisfying a $\delta$-Scaling inequality. Let $\gamma > 0$ and assume $x_0 \in K$ is such that $e^{-\gamma}w(x_0, \mathcal{S}_0)/2 \leq C_{f,s}^A$. With $\gamma_k = \gamma$, the output of Algorithm 2 satisfies*

$$f(x_T) - f^* \leq w_0 \frac{1}{\left(1 + \tilde{T}C_\gamma^\tau\right)^{\frac{1}{\tau}}} \quad \text{when } 1 \leq r < \frac{s}{s-1}$$

*after $T$ steps, with $w_0 \triangleq w(x_0, \mathcal{S}_0)$, $\tilde{T} \triangleq T - (|\mathcal{S}_0| - |\mathcal{S}_T|)$, and $\tau \triangleq \frac{s}{s-1} - r$. Also*

$$C_\gamma^\tau \triangleq \frac{e^{\gamma\tau} - 1}{C_s e^{\frac{s}{s-1}\gamma} w(x_0)^\tau},$$

*with $C_s \triangleq 2^{1+\frac{s}{s-1}} \frac{s}{s-1} \frac{c}{\delta} \big(C_{f,s}^A\big)^{\frac{1}{s-1}}$.*

13

*Proof.* Denote by $R$ the number of restarts after $T$ total inner iterations. We get

$$T \le \sum_{i=0}^{R-1} |\mathcal{S}_i| - |\mathcal{S}_{i+1}| + 2^{1+\frac{s}{s-1}} \frac{s}{s-1} e^{\frac{s}{s-1}\gamma} \big(C_{f,s}^A\big)^{\frac{1}{s-1}} \mu w(x_i, \mathcal{S}_i)^{r-\frac{s}{s-1}}.$$

Since $w(x_i, \mathcal{S}_i) \le w_0 e^{-\gamma i}$, it follows that

$$T \le |\mathcal{S}_0| - |\mathcal{S}_T| + 2^{1+\frac{s}{s-1}} \frac{s}{s-1} e^{\frac{s}{s-1}\gamma} \big(C_{f,s}^A\big)^{\frac{1}{s-1}} \mu w_0^{r-\frac{s}{s-1}} \sum_{i=0}^{R-1} e^{-\gamma i(r-\frac{s}{s-1})}.$$

Write $C_s = 2^{1+\frac{s}{s-1}} \frac{s}{s-1} \big(C_{f,s}^A\big)^{\frac{1}{s-1}} \mu$ and $\tau = \frac{s}{s-1} - r$ we have

$$T \le |\mathcal{S}_0| - |\mathcal{S}_T| + C_s e^{\frac{s}{s-1}\gamma} w_0^{r-\frac{s}{s-1}} \frac{e^{\gamma R \tau} - 1}{e^{\gamma \tau} - 1},$$

it follows that

$$e^{-\gamma R} \le \frac{1}{\left(1 + (T - (|\mathcal{S}_0| - |\mathcal{S}_T|)) \frac{(e^{\gamma \tau}-1)}{C_s e^{\frac{s}{s-1}\gamma} w(x_0)^\tau}\right)^{\frac{1}{\tau}}}.$$

which yields the desired result. ∎

Note that $r < \frac{s}{s-1}$ is always ensured because $s \in ]1, 2]$. In particular we only get linear convergence when $r = s = 2$ as for gradient methods [45]. We now show, as in Proposition 4.2, that the assumption $e^{-\gamma} w(x_0, \mathcal{S}_0)/2 \le C_f^A$ has a marginal impact on complexity when the function is $(L, s)$-Hölder smooth.

**Proposition 6.3** (Burn-in phase for Hölder smooth functions)**.** *After at most*

$$4 \frac{s}{s-1} \frac{e^\gamma}{\gamma} \ln \big(\frac{w_0}{2C_{f,s}^A}\big) + |\mathcal{S}_0|$$

*cumulative iterations of Algorithm 1, with constant schedule parameter $\gamma > 0$, we get a point $x$ such that $e^{-\gamma} w(x, \mathcal{S})/2 \le C_{f,s}^A$ when $f$ is $(L, s)$-Hölder smooth with $s > 1$.*

*Proof.* Assume we have $e^{-\gamma} w_0/2 > C_{f,s}^A$. Classically, the curvature argument ensures that we have for non-drop steps

$$\begin{aligned} f(x_t) - f(x_{t+1}) &\ge& \eta_t e^{-\gamma} w_0/2 - \frac{\eta_t^s}{s} C_{f,s}^A \\ &\ge& e^{-\gamma} w_0/2(1 - 1/s). \end{aligned}$$

Besides, $T_f$ being the number of full steps and $T$ the number of iterations before Fractional Away-step Frank-Wolfe stops,

$$f(x_0) - f(x_T) \ge T_f e^{-\gamma} w_0/2(1 - 1/s).$$

Combining this with $f(x_0) - f(x_T) \le f(x_0) - f(x^*) \le w_0$ we get

$$T_f \le 2e^\gamma \frac{s}{s-1}.$$

Finally with the classical counting argument on drop steps, we obtain

$$T \le 4e^\gamma \frac{s}{s-1} + |\mathcal{S}_0| - |\mathcal{S}_T|.$$

Denote $R$ the number of calls to the Fractional Away-step Frank-Wolfe before the last output $\hat{x}_R$ satisfies $e^{-\gamma} w(\hat{x}, \mathcal{S}_{\hat{x}})/2 > C_{f,s}^A$. The strong Wolfe gap of the $N^{th}$ output of the Fractional Away-step Frank-Wolfe satisfies by definition

$$w(\hat{x}_N) \le e^{-N\gamma} w_0,$$

hence we have

$$R \le \frac{1}{\gamma} \ln \big(\frac{w_0}{2C_{f,s}^A}\big).$$

14

Finally each round of the Fractional Away-step Frank-Wolfe under the initial assumption that $e^{-\gamma}w(\hat{x}_i, \mathcal{S}_{\hat{x}_i})/2 > C_{f,s}^A$ require at most $4e^\gamma \frac{s}{s-1} + |\mathcal{S}_{\hat{x}_i}| - |\mathcal{S}_{\hat{x}_{i+1}}|$ iterations. Hence a total $T_t$ of

$$
\begin{aligned}
T_t &\leq \sum_{i=1}^R 4e^\gamma \frac{s}{s-1} + |\mathcal{S}_{\hat{x}_i}| - |\mathcal{S}_{\hat{x}_{i+1}}| \\
&\leq 4Re^\gamma \frac{s}{s-1} + |\mathcal{S}_0| \\
&\leq 4\frac{s}{s-1}\frac{e^\gamma}{\gamma} \ln\left(\frac{w_0}{2C_{f,s}^A}\right) + |\mathcal{S}_0|
\end{aligned}
$$

which is the desired result. ∎

## 7. Fractional Frank-Wolfe Algorithm

In this section, we describe how Hölderian error bounds coupled with a restart scheme yield improved convergence bounds for the vanilla Frank-Wolfe algorithm.

In Sections 4-6, relaxing strong convexity of $f$ using the $(\theta, c)$-HEB assumption lead to improved sub-linear rates using a restart scheme for the Away step variant of the Frank-Wolfe algorithm when the set of constraints $\mathcal{C}$ is a polytope. For these sets, away steps produce accelerated convergence rates that the vanilla Frank-Wolfe algorithm cannot achieve.

However, accelerated convergence hold for the vanilla Frank-Wolfe algorithm in other scenarios. For instance, when the solution of (2) is in the interior of the set and $f$ is strongly convex, the convergence of the vanilla Frank-Wolfe is linear. In this vein, we define a fractional version of the Frank-Wolfe algorithm (Algorithm 3) and analyse its restart scheme (Algorithm 4) under the $(\theta, c)$-HEB condition in Section 7.2. Although the Fractional Frank-Wolfe algorithm and the vanilla Frank-Wolfe algorithm perform the same iterations, the restart scheme produces a much simpler proof of improved convergence bounds. The fractional variant is also the structural basis for recent competitive versions of the Frank-Wolfe algorithm [7].

Another acceleration scenario for the vanilla Frank-Wolfe algorithm is when the set of constraints $\mathcal{C}$ is strongly convex. Under some restrictive assumption on $f$, the classical analysis [34, (5) in Theorem 6.1] exhibits a linear convergence rate. Recently [18] have shown a general $\mathcal{O}(1/T^2)$ sub-linear rate when $f$ and $\mathcal{C}$ are strongly convex. We will state new rates for the case where $f$ satisfies $(\theta, c)$-HEB and $\mathcal{C}$ is strongly convex, to provide a complete picture.

For completeness, we would like to mention that $\delta$-scaling for the away step Frank-Wolfe algorithm does not apply in the case where $\mathcal{C}$ is a strongly convex set. Lemma 3.4 does not hold anymore, and $PWidth$ can tend to zero in this case.

**7.1. Restart schemes for Fractional Frank-Wolfe.** We now state the *fractional* version of the (vanilla) Frank-Wolfe algorithm. The Fractional Frank-Wolfe algorithm 3 is derived from Algorithm 1 by replacing $w(x_0, \mathcal{S}_0)$ with $g(x_0)$, as in (6) and dropping the away step update.

---
**Algorithm 3** Fractional Frank-Wolfe Algorithm

**Input:** A smooth convex function $f$. Starting point $x_0 \in \mathcal{C}$. LP oracle (3) and schedule parameter $\gamma > 0$.
1: $t := 0$
2: **while** $g(x_t) > e^{-\gamma}g(x_0)$ **do**
3: $\quad v_t := \mathrm{LP}_\mathcal{C}(\nabla f(x_t))$ and $d_t^{FW} := v_t - x_t$
4: $\quad x_{t+1} := x_t + \eta_t d_t^{FW}$ with $\eta_t \in [0,1]$ via line-search
5: $\quad t := t + 1$
6: **end while**
**Output:** $x_t \in \mathcal{C}$ such that $g(x_t) \leq e^{-\gamma}g(x_0)$

---

A constant restart scheme using Algorithm 3 for its inner iteration, recovers the Scaling Frank-Wolfe algorithm [7, Algorithm 7: Parameter-free Lazy Conditional Gradient] up to a slight reformulation with the additional $\Phi_t$ parameter. The two algorithms have the same restart structure. However, the Scaling Frank-Wolfe algorithm additionally uses a weaker oracle (a so-called Weak Separation Oracle) than the Linear Optimization Oracle that we employ here. More precisely, the Scaling Frank-Wolfe algorithm does not necessarily require $v_t$ to be the exact nor an approximate solution of the Linear Minimization Problem, but rather to satisfy the condition $\langle -\nabla f(x_t), v_t - x_t \rangle > \Phi_t e^{-\gamma}$. As a consequence, $g(x_t)$ is not computed and $\Phi_t$ is only an upper bound on $g(x_t)$. This explains the difference in line 8 of Algorithm 4.

---

**Algorithm 4** Restart Fractional Frank-Wolfe Algorithm

---

**Input:** A smooth convex function $f$ with curvature $C_f$. Starting point $x_0 \in \mathcal{C}$. $\epsilon > 0$, LP oracle (3) and schedule parameter $\gamma > 0$.

1: $t := 0$ and $\Phi_0 := g(x_0)$
2: **while** $g(x_t) > \epsilon$ **do**
3: $\quad v_t := \mathrm{LP}_{\mathcal{C}}(\nabla f(x_t))$ and $d_t^{FW} := v_t - x_t$
4: $\quad$ **if** $\langle -\nabla f(x_t), v_t - x_t \rangle > \Phi_t e^{-\gamma}$ **then**
5: $\qquad x_{t+1} := x_t + \eta_t d_t^{FW}$ with $\eta_t \in [0, 1]$ via line-search
6: $\qquad \Phi_{t+1} := \Phi_t$
7: $\quad$ **else**
8: $\qquad \Phi_{t+1} := g(x_t)$ (hence $\Phi_{t+1} < \Phi_t e^{-\gamma}$)
9: $\quad$ **end if**
10: $\quad t := t + 1$
11: **end while**

---

### 7.2. Optimum in the Interior of the Feasible Set.

We first recall that when the optimal solutions of (1) are in the interior of $\mathcal{C}$, a version of the (Scaling) inequality is automatically satisfied. (FW-Scaling) replaces $w(x)$ by $g(x)$ and can be interpreted as a scaling inequality tailored to the (vanilla) Frank-Wolfe algorithm. Note that the $\delta$ parameter depends on the relative distance of the optimal set $X^*$ to the boundary of $\mathcal{C}$. This property has already been extensively used in, e.g., [22, 19, 20].

**Lemma 7.1** (FW $\delta$-scaling when optimum is in interior [22]). *Assume $\mathcal{C}$ is convex and $f$ convex differentiable. Assume $X^* \subset \mathbf{Int}(\mathcal{C})$ and choose $z > 0$ such that $B(x^*, z) \subset \mathcal{C}$ for all $x^* \in X^*$. Then for all $x \in \mathcal{C}$ such that $d(x, X^*) \leq \frac{z}{2}$ we have*

$$g(x) \geq \frac{z}{2} \|\nabla f(x)\|, \qquad \text{(FW-Scaling)}$$

*where $g(x)$ is the Frank-Wolfe (dual) gap as defined in (6).*

*Proof.* For $x \in B(x^*, \frac{z}{2})$, we have $x - \frac{z}{2} \frac{\nabla f(x)}{\|\nabla f(x)\|} \in \mathcal{C}$. Denote $v$ the Frank-Wolfe vertex, we have $g(x) \triangleq \langle -\nabla f(x), v - x \rangle$. By optimality of $v$, we hence obtain

$$g(x) \geq \langle -\nabla f(x), x - \frac{z}{2} \frac{\nabla f(x)}{\|\nabla f(x)\|} - x \rangle = \frac{z}{2} \|\nabla f(x)\|,$$

which is the desired result. ∎

We now bound the convergence rate of Algorithm 4 in the following proposition.

**Proposition 7.2** (Convergence Rate of Restart Fractional FW). *Let $f$ be a smooth convex function with curvature $C_f$ as defined in (10), satisfying $(\theta, c)$-HEB on $\mathcal{C}$. Assume there exists $z > 0$ such that $B(x^*, z) \subset \mathcal{C}$ for all $x^* \in X^*$. Let $\gamma > 0$ and assume $x_0$ is such that $e^{-\gamma} g(x_0) \leq C_f$ and $f(x_0) - f^* \leq \left(\frac{z}{2}\right)^{\frac{1}{\theta}}$ (burn-in*

16

*phase). Then the output of Algorithm 2 satisfies ($r = \frac{1}{1-\theta}$)*

$$\begin{cases} f(x_T) - f^* \leq g_0 \dfrac{1}{\left(1 + TC_\gamma^r\right)^{\frac{1}{2-r}}} & \text{when } 1 \leq r < 2 \\[4mm] f(x_T) - f^* \leq g_0 \exp\left(-\dfrac{\gamma}{e^{2\gamma}} \dfrac{T}{8C_f\mu}\right) & \text{when } r = 2 \ , \end{cases}$$

*after $T$ steps, with $g_0 = g(x_0)$. Also, with $\mu = (cz/2)^{1/(1-\theta)}$ we write*

$$C_\gamma^r \triangleq \frac{e^{\gamma(2-r)} - 1}{2e^{2\gamma}C_f\mu g(x_0)^{r-2}}.$$

*Proof.* First note that for all $t$, we have $d(x_t, X^*) \leq \frac{z}{2}$. Indeed $f(x_t) - f^* \leq f(x_{t-1}) - f^* \leq \left(\frac{z}{2}\right)^{\frac{1}{\theta}}$. Hence by $(\theta, c)$-HEB we have

$$\min_{x^* \in X^*} \|x_t - x^*\| \leq (f(x_t) - f^*)^\theta \leq \frac{z}{2}.$$

We can now apply lemma 7.1 to get for all $x_t$

$$g(x_t) \geq \frac{z}{2}\|\nabla f(x_t)\|,$$

and as in Lemma 3.6, FW-Scaling and $(\theta, c)$-HEB leads to a Wolfe primal gap (with $\mu = (cz/2)^{1/(1-\theta)} > 0$)

$$f(x) - f^* \leq \mu g(x)^r,$$

where $r = 1/(1 - \theta)$. The proof then follows exactly that of Fractional Away Frank-Wolfe and its restart schemes (see Proposition 4.1 and Theorem 5.1), replacing $w(x)$ with $g(x)$. The only change comes from the upper bound on $T$, the number of iterations needed for Fractional Frank-Wolfe to stop. We recall the key steps to get this bound and update its value. At each iteration

$$f(x_t) - f(x_{t+1}) \geq \max_{\eta \in [0,1]} \{\eta e^{-\gamma} g(x_0) - \frac{\eta^2}{2}C_f\},$$

such that because of assumption $e^{-\gamma}g(x_0) < C_f$, we have

$$f(x_t) - f(x_{t+1}) \geq \frac{1}{2} \frac{g(x_0)^2}{e^{2\gamma}C_f}.$$

Hence on one side

$$f(x_0) - f(x_T) \geq \frac{T}{2} \frac{g(x_0)^2}{e^{2\gamma}C_f}.$$

And on the other side, using the $r$-Wolfe primal bound $f(x_0) - f(x_T) \leq \mu g(x_0)^r$ and finally

$$T \leq 2\mu C_f e^{2\gamma} g(x_0)^{r-2}.$$

The restart scheme is then controlled exactly as in the proof of 5.1. ∎

Assuming that $e^{-\gamma}g(x_0) \leq C_f$ and $f(x_0) - f^* \leq \left(\frac{z}{2}\right)^{\frac{1}{\theta}}$ simplify the statements and it is automatically satisfied after a burn-in phase. However it is fundamental to assume that there exists $z > 0$ s.t. $B(x^*, z) \subset C$ for all $x^* \in X^*$. Indeed this ensures that the optimal set is in the interior of $C$. Note also that a robustness result similar to that of Proposition 5.3 holds here.

**7.3. Strongly Convex Constraint Set.** When $\mathcal{C}$ is strongly convex, strong convexity of $f$ leads to a better convergence rate than the sub-linear $\mathcal{O}(1/T)$. The original analysis of [34, (5) in Theorem 6.1] assumes $\|\nabla f(x)\| \geq \epsilon > 0$ (irrespective of the strong convexity of $f$) and hence $(\theta, c)$-HEB cannot be understood as a relaxation of the assumption. This analysis provides a linear convergence rate when the unconstrained minimum of $f$ is strictly outside of $\mathcal{C}$. §7.2 shows linear convergence when $x^*$ is in the interior of $\mathcal{C}$. Hence the remaining case is when the unconstrained minimum of $f$ is in $\partial \mathcal{C}$, the boundary of $\mathcal{C}$ (an arguably rare instance).

Recently, the analysis of [18] closes this gap by providing a general convergence rate of $\mathcal{O}(1/T^2)$ under a (slightly) weaker assumption than strong convexity of $f$ [18, see (2)]. Although the asymptotic rate regime of [18] is significantly less appealing than the linear convergence rate in [34] and hence seemingly a marginal improvement in term of applicability, the situation is a bit more complicated: the bound of [18] benefits from much better conditioning and can easily dominate other bounds near $\partial \mathcal{C}$. In particular, the conditioning of [34] depends on the $\epsilon$ lower bounding the norm of the gradient on the constraint set, which can be arbitrarily small. The analysis of [18] adapts to $(\theta, c)$-HEB, as was detailed in [47] and we recall this below for the sake of completeness.

**Theorem 7.3.** *Consider $\mathcal{C}$ an $\alpha$-strongly convex set and $f$ a convex $L$-smooth function* (8). *Assume $(\theta, c)$-HEB for $f$. Then the iterate of (vanilla) Frank-Wolfe is such that $f(x_T) - f(x^*) = \mathcal{O}\left(1/T^{1/(1-\theta)}\right)$ for $\theta \in [0, 1/2]$.*

*Proof.* From [18, Lemma 1], $L$-smoothness of $f$ combines with $\alpha$-strong convexity of $\mathcal{C}$ gives

$$h_{t+1} \leq h_t \cdot \max\left\{\frac{1}{2}, 1 - \frac{\alpha\|\nabla f(x)\|}{8L}\right\}.$$

On the other hand with $(\theta, c)$-HEB and by convexity of $f$, (12) applies

$$\left(f(x) - f(x^*)\right)^{1-\theta} \leq c \cdot \min_{y \in X^*} \frac{\langle \nabla f(x), \, x - x^* \rangle}{\|x - x^*\|}$$
$$\leq c \, \|\nabla f(x)\|.$$

Note that with $\theta = 1/2$, this is the sufficient condition [18, (2)] implied by strong convexity that leads to $\mathcal{O}(1/T^2)$ convergence rates. Hence combining both we recover this recursive inequality for $h_t = f(x_t) - f(x^*)$

$$h_{t+1} \leq h_t \cdot \max\left\{\frac{1}{2}, 1 - \frac{\alpha}{8Lc} h_t^{1-\theta}\right\}.$$

When $\theta = 0$ (convexity), this leads to the classical $\mathcal{O}(1/T)$ rate. When $\theta = 1/2$ the above recursion leads to a $\mathcal{O}(1/T^2)$ rate as in [18, proof of Theorem 2]. With Lemma B.1 in Appendix B, for any non-negative constants $(k, C)$, such that $\frac{2 - 2^\beta}{2^\beta - 1} \leq k$ and $\max\{h_0 k^{1/\beta}, \frac{2}{\left((\beta - (1-\beta)(2^\beta - 1))M\right)^{1/\beta}}\} \leq C$ (with $M \triangleq \frac{\alpha}{8L}$), we have

$$h_t \leq \frac{C}{(t+1)^{1/(1-\theta)}}.$$

and the desired result. ∎

Theorem 7.3 interpolates between the general $\mathcal{O}(1/T)$ rate for smooth convex functions and the $\mathcal{O}(1/T^2)$ rate for smooth and strongly convex functions.

## 8. CONCLUSION

We derived a variant of the Away-step Frank-Wolfe algorithm and showed improved complexity bounds when the strong Wolfe gap satisfies a generalized strong convexity condition. The Łojasiewicz factorization lemma shows that this condition actually holds generically for some value of the parameters, producing complexity bounds of the form $O(1/\epsilon^q)$ with $q \leq 1$, thus smoothly interpolating between the complexity of the classical FW algorithm with rate $O(1/\epsilon)$ and that of the Away-step Frank-Wolfe with rate $O(\log(1/\varepsilon))$.

Our method is adaptive to the value of the generalized strong convexity parameters and robustly yields optimal performance.

## References

[1] H. Attouch, J. Bolte, P. Redont, and A. 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.

[2] H. Attouch, G. Buttazzo, and G. Michaille. Variational analysis in Sobolev and BV spaces: applications to PDEs and optimization. *SIAM*, 17, 2014.

[3] M. A. Bashiri and X. Zhang. Decomposition-invariant conditional gradient for general polytopes with line search. In *Advances in Neural Information Processing Systems*, pages 2687–2697, 2017.

[4] H. H. Bauschke, J. Bolte, and M. Teboulle. A descent lemma beyond lipschitz gradient continuity: first-order methods revisited and applications. *Mathematics of Operations Research*, 42(2):330–348, 2016.

[5] J. Bolte, A. Daniilidis, and A. Lewis. The Łojasiewicz inequality for nonsmooth subanalytic functions with applications to subgradient dynamical systems. *SIAM Journal on Optimization*, 17(4):1205–1223, 2007.

[6] J. Bolte, T. P. Nguyen, J. Peypouquet, and B. W. Suter. From error bounds to the complexity of first-order descent methods for convex functions. *Mathematical Programming*, 165(2):471–507, 2017.

[7] G. Braun, S. Pokutta, and D. Zink. Lazifying conditional gradient algorithms. In *International Conference on Cachine Learning*, pages 566–575, 2017.

[8] G. Braun, S. Pokutta, D. Tu, and S. Wright. Blended conditonal gradients. In *International Conference on Machine Learning*, pages 735–743, 2019.

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

[10] Z. Chen, Y. Xu, E. Chen, and T. Yang. Sadagrad: Strongly adaptive stochastic gradient methods. In *International Conference on Machine Learning*, pages 912–920, 2018.

[11] C. Combettes and S. Pokutta. Boosting Frank-Wolfe by chasing gradients. In *International Conference on Machine Learning*, pages 2111–2121, 2020.

[12] J. Diakonikolas, A. Carderera, and S. Pokutta. Locally accelerated conditional gradients. In *International Conference on Artificial Intelligence and Statistics*, pages 1737–1747, 2020.

[13] O. Fercoq and Z. Qu. Restarting accelerated gradient methods with a rough strong convexity estimate. *arXiv preprint arXiv:1609.07358*, 2016.

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

[15] P. Frankel, G. Garrigos, and J. Peypouquet. Splitting methods with variable metric for Kurdyka-Łojasiewicz functions and general convergence rates. *Journal of Optimization Theory and Applications*, 165(3):874–900, 2015.

[16] R. M. Freund and P. Grigas. New analysis and results for the Frank-Wolfe method. *Mathematical Programming*, 155(1):199–230, 2016. ISSN 1436-4646.

[17] R. M. Freund, P. Grigas, and R. 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.

[18] D. Garber and E. Hazan. Faster rates for the Frank-Wolfe method over strongly-convex sets. In *Proceedings of the 32th International Conference on Machine Learning*, 2015.

[19] D. Garber and E. Hazan. A linearly convergent variant of the conditional gradient algorithm under strong convexity, with applications to online and stochastic optimization. *SIAM Journal on Optimization*, 26(3):1493–1528, 2016.

[20] D. Garber and O. Meshi. Linear-memory and decomposition-invariant linearly convergent conditional gradient algorithm for structured polytopes. In *Advances in Neural Information Processing Systems*, 2016.

[21] P. Giselsson and S. Boyd. Monotonicity and restart in fast gradient methods. In *53rd IEEE Conference on Decision and Control*, pages 5058–5063. IEEE, 2014.

[22] J. Guélat and P. Marcotte. Some comments on Wolfe's 'away step'. *Mathematical Programming*, 35(1):110–119, 1986.

[23] A. J. Hoffman. On approximate solutions of systems of linear inequalities. *Journal of Research of the National Bureau of Standards*, 49(4), 1952.

[24] M. Jaggi. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In *Proceedings of the 30th International Conference on Machine Learning*, pages 427–435, 2013.

[25] A. Joulin, K. Tang, and L. Fei-Fei. Efficient image and video co-localization with Frank-Wolfe algorithm. In *European Conference on Computer Vision*, pages 253–268. Springer, 2014.

[26] H. Karimi, J. Nutini, and M. 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.

[27] T. Kerdreux, F. Pedregosa, and A. d'Aspremont. Frank-Wolfe with subsampling oracle. In *International Conference on Machine Learning*, pages 2591–2600, 2018.

[28] T. Kerdreux, A. d'Aspremont, and S. Pokutta. Restarting Frank-Wolfe. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1275–1283, 2019.

[29] T. Kerdreux, A. d'Aspremont, and S. Pokutta. Projection-free optimization on uniformly convex sets. In *International Conference on Artificial Intelligence and Statistics*, pages 19–27, 2021.

[30] S. Lacoste-Julien and M. Jaggi. On the global linear convergence of Frank-Wolfe optimization variants. In *Advances in Neural Information Processing Systems*, volume 28, pages 496–504, 2015.

[31] S. Lacoste-Julien, M. Jaggi, M. Schmidt, and P. Pletscher. Block-coordinate Frank-Wolfe optimization for structural SVMs. In *International Conference on Machine Learning*, pages 53–61, 2013.

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

[33] G. Lan, S. Pokutta, Y. Zhou, and D. Zink. Conditional accelerated lazy stochastic gradient descent. In *International Conference on Machine Learning*, pages 1965–1974, 2017.

[34] E. S. Levitin and B. T. Polyak. Constrained minimization methods. *USSR Computational Mathematics and Mathematical Physics*, 6(5):1–50, 1966.

[35] G. Li and T. K. 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.

[36] F. Locatello, R. Khanna, M. Tschannen, and M. Jaggi. A unified optimization view on generalized matching pursuit and Frank-Wolfe. In *Artificial Intelligence and Statistics*, pages 860–868, 2017.

[37] F. Locatello, M. Tschannen, G. Rätsch, and M. Jaggi. Greedy algorithms for cone constrained optimization with convergence guarantees. In *Advances in Neural Information Processing Systems*, pages 773–784, 2017.

[38] S. Lojasiewicz. Ensembles semi-analytiques. *Institut des Hautes Études Scientifiques*, 1965.

[39] S. Lojasiewicz. Sur la géométrie semi-et sous-analytique. *Ann. Inst. Fourier*, 43(5):1575–1595, 1993.

[40] A. Miech, J.-B. Alayrac, P. Bojanowski, I. Laptev, and J. Sivic. Learning from video and text via large-scale discriminative clustering. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 5267–5276. IEEE, 2017.

[41] A. Nemirovskii and Y. E. Nesterov. Optimal methods of smooth convex minimization. *USSR Computational Mathematics and Mathematical Physics*, 25(2):21–30, 1985.

[42] B. O'Donoghue and E. Candes. Adaptive restart for accelerated gradient schemes. *Foundations of Computational Mathematics*, 15(3):715–732, 2015.

[43] A. Osokin, J.-B. Alayrac, I. Lukasewitz, P. K. Dokania, and S. Lacoste-Julien. Minding the gaps for block Frank-Wolfe optimization of structured SVMs. *International Conference on Machine Learning*, 2016.

[44] J. Pena and D. Rodriguez. Polytope conditioning and linear convergence of the Frank-Wolfe algorithm. *Mathematics of Operations Research*, 2018.

[45] V. Roulet and A. d'Aspremont. Sharpness, restart, and acceleration. *SIAM Journal on Optimization*, 30(1): 262–289, 2020.

[46] N. Shah, V. Kolmogorov, and C. H. Lampert. A multi-plane block-coordinate Frank-Wolfe algorithm for training structural SVMs with a costly max-oracle. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2737–2745, 2015.

[47] Y. Xu and T. Yang. Frank-Wolfe method is automatically adaptive to error bound condition. *arXiv:1810.04765*, 2018.

## Appendix A. One shot application of the Fractional Away-step Frank Wolfe

Fractional Away-step Frank-Wolfe output a point $x_t \in \mathcal{C}$ s.t. $w(x_t) \leq e^{-\gamma} w_0$. Hence, running once Fractional Away-step Frank-Wolfe with a large value of $\gamma$ allows finding an approximate minimizer with the desired precision. The following lemma proves a sublinear $\mathcal{O}(1/T)$ convergence rate which corresponds to the convergence rate of the Frank-Wolfe algorithm. Importantly the rate does not depend on $r$. Hence there is no hope of observing linear convergence for the strongly convex case.

**Lemma A.1.** *Let $f$ be a smooth convex function, $\epsilon > 0$ be a target accuracy, and $x_0 \in \mathcal{C}$ be an initial point. Then for any $\gamma > \ln \frac{w(x_0)}{\epsilon}$, Algorithm 1 satisfies:*

$$f(x_T) - f^* \leq \epsilon,$$

*for $T \geq \frac{2C_f^{\mathcal{A}}}{\epsilon}$.*

*Proof.* We can stop the algorithm as soon as the criterion $w(x_t) < \epsilon$ in step 2 is met or we observe an away step, whichever comes first. In former case we have $f(x_t) - f^* \leq w(t) < \epsilon$, in the latter it holds

$$f(x_t) - f^* \leq -\nabla f(x_t)(d_t^{FW}) \leq \epsilon/2 < \epsilon.$$

Thus, when the algorithms stops, we have achieved the target accuracy and it suffices to bound the number of iterations required to achieve that accuracy. Moreover, while running, the algorithm only executes Frank-Wolfe and we drop the FW superscript in the directions; otherwise we would have stopped.

From the proof of Proposition 4.1, we have each Frank-Wolfe step ensures progress of the form

$$f(x_t) - f(x_{t+1}) \geq \begin{cases} \frac{\langle -\nabla f(x_t), d_t \rangle^2}{2C_f^{\mathcal{A}}} & \text{if } \langle -\nabla f(x_t), d_t \rangle \leq C_f^{\mathcal{A}} \\ \langle -\nabla f(x_t), d_t \rangle - C_f^{\mathcal{A}}/2 & \text{otherwise.} \end{cases}$$

For convenience, let $h_t \triangleq f(x_t) - f^*$. By convexity we have $h_t \leq \langle -\nabla f(x_t), d_t \rangle$, so that the above becomes

$$f(x_t) - f(x_{t+1}) \geq \begin{cases} \frac{h_t^2}{2C_f^{\mathcal{A}}} & \text{if } h_t \leq C_f^{\mathcal{A}} \\ h_t - C_f^{\mathcal{A}}/2 & \text{otherwise.} \end{cases},$$

and moreover observe that the second case can only happen in the very first step: $h_1 \leq h_0 - (h_0 - C_f^{\mathcal{A}}/2) = C_f^{\mathcal{A}}/2 \leq 2C_f^{\mathcal{A}}/t$ for $t = 1$ providing the start of the following induction: we claim $h_t \leq \frac{2C_f^{\mathcal{A}}}{t}$.

Suppose we have established the bound for $t$, then for $t + 1$, we have

$$h_{t+1} \leq \left(1 - \frac{h_t}{2C_f^{\mathcal{A}}}\right) h_t \leq \frac{2C_f^{\mathcal{A}}}{t} - \frac{2C_f^{\mathcal{A}}}{t^2} \leq \frac{2C_f^{\mathcal{A}}}{t+1}.$$

The induction is complete and it follows that the algorithm requires $T \geq \frac{2C_f^{\mathcal{A}}}{\epsilon}$ to reach $\epsilon$-accuracy. ∎

## Appendix B. Non Standard Recurrence Relation

This is a technical Lemma derived in [47, proof of Theorem 1. ] and we repeat it here for the sake of completeness. It is used in Section 7.

**Lemma B.1** (Recurrence and sub-linear rates)**.** *Consider a sequence $(h_n)$ of non-negative numbers. Assume there exists $M > 0$ and $0 < \beta \leq 1$ s.t.*

$$h_{t+1} \leq h_t \max\{1/2, 1 - Mh_t^\beta\}, \tag{20}$$

*then $h_T = \mathcal{O}(1/T^{1/\beta})$. More precisely for all $t \geq 0$,*

$$h_t \leq \frac{C}{(t+k)^{1/\beta}}$$

*with $(k, C)$ such that $\frac{2-2^\beta}{2^\beta-1} \leq k$ and $\max\{h_0 k^{1/\beta}, 2(C'/M)^{1/\beta}\} \leq C$, with $C' \geq \frac{1}{\beta-(1-\beta)(2^\beta-1)}$.*

*Proof.* Let $(k, C)$ satisfying the condition in Lemma B.1. Let's show by induction that

$$h_t \leq \frac{C}{(t+k)^{1/\beta}} \ .$$

For $t = 0$, it is true because we assumed $h_0 k^{1/\beta} \leq C$. Let $t \geq 1$. Consider the case where the maximum in the right hand side of (20) is obtained with $\frac{1}{2}$, then

$$
\begin{aligned}
h_{t+1} &\leq \frac{C}{(t+k+1)^{1/\beta}}\Big(\frac{t+k+1}{t+k}\Big)^{1/\beta}\frac{1}{2} \\
&\leq \frac{C}{(t+k+1)^{1/\beta}},
\end{aligned}
$$

because $k \geq \frac{2-2^\beta}{2^\beta-1}$. Otherwise we have

$$h_{t+1} \leq h_t(1 - M h_t^\beta).$$

If $h_t \leq \frac{C}{2(t+k)^{1/\beta}}$, conclusion holds as before. Otherwise assume $\frac{C}{2(t+k)^{1/\beta}} \leq h_t \leq \frac{C}{(t+k)^{1/\beta}}$ and (20) implies

$$
\begin{aligned}
h_{t+1} &\leq \frac{C}{(t+k)^{1/\beta}}\Big(1 - M\Big(\frac{C}{2}\Big)^\beta \frac{1}{t+k}\Big) \\
h_{t+1} &\leq \frac{C}{(t+k+1)^{1/\beta}}\Big(1 + \frac{1}{t+k}\Big)^{1/\beta}\Big(1 - M\Big(\frac{C}{2}\Big)^\beta \frac{1}{t+k}\Big)
\end{aligned}
$$

From Lemma B.2, for $C \geq 2(C'/M)^{1/\beta}$ with $C' \geq \frac{1}{\beta-(1-\beta)(2^\beta-1)}$, we have

$$h_{t+1} \leq \frac{C}{(t+k+1)^{1/\beta}}\Big(1 + \frac{C'}{t+k}\Big)\Big(1 - \frac{C'}{t+k}\Big) \leq \frac{C}{(t+k+1)^{1/\beta}} \ ,$$

which proves the induction. ∎

**Lemma B.2.** *For any $t \geq 1$, we have*

$$\Big(1 + \frac{1}{t+k}\Big)^{1/\beta}\Big(1 - M\Big(\frac{C}{2}\Big)^\beta \frac{1}{t+k}\Big) \leq \Big(1 + \frac{C'}{t+k}\Big)\Big(1 - \frac{C'}{t+k}\Big), \tag{21}$$

*where $k \geq \frac{2-2^\beta}{2^\beta-1}$, $\beta \in ]0, 1]$, $C' \geq \frac{1}{\beta-(1-\beta)(2^\beta-1)}$ and $C$ such that $C \geq 2(C'/M)^{1/\beta}$.*

*Proof.* Write $x = \frac{1}{t+k}$. Because $t \geq 1$ and $k \geq \frac{2-2^\beta}{2^\beta-1}$, we have $x \in ]0, 2^\beta - 1]$. (21) is equivalent to

$$\frac{1}{\beta}\log(1 + \frac{1}{t+k}) + \log(1 - M\Big(\frac{C}{2}\Big)^\beta x) \leq \log(1 + C'x) + \log(1 - C'x)$$

Choosing $C$ greater or equal to $2\Big(C'/M\Big)^{1/\beta}$ ensures that $\log(1 - M\Big(\frac{C}{2}\Big)^\beta x) \leq \log(1 - C'x)$. Also for $C' \geq \frac{1}{\beta-(1-\beta)(2^\beta-1)}$, the function $h(x) \triangleq \log(1 + C'x) - \frac{1}{\beta}\log(1 + x)$ is non-decreasing (and hence non-negative) on $]0, 2^\beta - 1]$. Reciprocally for $C' \geq \frac{1}{\beta-(1-\beta)(2^\beta-1)}$ and $C \geq \Big(C'/M\Big)^{1/\beta}$, (21) holds. ∎

ZUSE INSTITUTE BERLIN & TECHNISCHE UNIVERSITÄT BERLIN, GERMANY
*Email address*: thomaskerdreux@gmail.com

CNRS & D.I., UMR 8548,
ÉCOLE NORMALE SUPÉRIEURE, PARIS, FRANCE.
*Email address*: aspremon@ens.fr

ZUSE INSTITUTE BERLIN & TECHNISCHE UNIVERSITÄT BERLIN, GERMANY
*Email address*: pokutta@zib.de