

---

# Affine Invariant Analysis of Frank-Wolfe on Strongly Convex Sets

---

Thomas Kerdreux<sup>1,\*</sup>

Lewis Liu<sup>2,\*</sup>

Simon Lacoste-Julien<sup>2,3,†</sup>

Damien Scieur<sup>3,\*</sup>

## Abstract

It is known that the Frank-Wolfe (FW) algorithm, which is affine-covariant, enjoys accelerated convergence rates when the constraint set is strongly convex. However, these results rely on norm-dependent assumptions, usually incurring non-affine invariant bounds, in contradiction with FW’s affine-covariant property. In this work, we introduce new structural assumptions on the problem (such as the directional smoothness) and derive an affine invariant, norm-independent analysis of Frank-Wolfe. Based on our analysis, we propose an affine invariant backtracking line-search. Interestingly, we show that typical backtracking line-searches using smoothness of the objective function surprisingly converge to an affine invariant step size, despite using affine-dependent norms in the step size’s computation. This indicates that we do not necessarily need to know the set’s structure in advance to enjoy the affine-invariant accelerated rate.

## 1 Introduction

Conditional Gradient algorithms, a.k.a. Frank-Wolfe (FW) algorithms [Frank et al. 1956], form a class of first-order methods solving constrained optimization problems such as

$$\min_{x \in \mathcal{C}} f(x). \quad (1)$$

The schemes in this class decompose non-linear constrained problems into a series of linear problems on the original constraint set, *i.e.* linear minimization oracles (LMO). They form a practical family of algorithms [Jaggi 2013; Bojanowski et al. 2014; Alayrac et al. 2016; Seguin et al. 2016; Peyre et al. 2017; Miech

---

## Algorithm 1 Frank-Wolfe Algorithm

---

**Input:**  $x_0 \in \mathcal{C}$ .

- 1: **for**  $k = 0, 1, \dots, K$  **do**
  - 2:  $v_k \in \operatorname{argmax}_{v \in \mathcal{C}} \langle -\nabla f(x_k), v - x_k \rangle \quad \triangleright$  LMO
  - 3:  $\gamma_k = \operatorname{argmin}_{\gamma \in [0,1]} f(x_k + \gamma(v_k - x_k)) \quad \triangleright$  Line-search
  - 4:  $x_{k+1} = (1 - \gamma_k)x_k + \gamma_k v_k \quad \triangleright$  Convex update
  - 5: **end for**
- 

et al. 2018; Lacoste-Julien et al. 2015b; Courty et al. 2016; Paty et al. 2019; Luise et al. 2019]; however, many open questions remain in designing such optimal algorithmic schemes (*e.g.* [Braun et al. 2017; Kerdreux et al. 2018b; Braun et al. 2018; Cyrille W Combettes et al. 2020b; Carderera et al. 2020; Mortagy et al. 2020; Cyrille W. Combettes et al. 2020a]) and in their theoretical understanding.

Besides, with the appropriate line-search, the iterates of the FW are *affine covariant* under the affine transformation  $y = Bx + b$  of problem (1),

$$\min_{y \in \tilde{\mathcal{C}} = B^{-1}(\mathcal{C} - b)} \tilde{f}(y) \stackrel{\text{def}}{=} f(B^{-1}(y - b)), \quad B \text{ invertible.} \quad (2)$$

**Definition 1.1** *An algorithm is affine covariant when its iterates  $(x_k)$  (resp.  $(y_k)$ ) for problem (1) (resp. (2)) satisfy*

$$y_k = Bx_k + b.$$

In other words, the behavior of Algorithm 1 is insensitive to affine transformations or re-parametrization of the space. This means that, ideally, the theoretical rate for an affine covariant algorithm should be *affine invariant*.

The original Frank-Wolfe algorithm (Algorithm 1) generally enjoy a slow sublinear rate  $\mathcal{O}(1/K)$  over general compact convex set and smooth convex functions [Jaggi 2013]. In that setting, Clarkson 2010; Jaggi 2013 define a modulus of smoothness that leads to affine invariant analysis of the Frank-Wolfe algorithm, matching with the affine covariant behavior of the algorithm.

---

\* Equal contribution

<sup>1</sup> Zuse Institute, Berlin.

<sup>2</sup> MILA and DIRO, Université de Montréal

<sup>3</sup> Samsung SAIT AI Lab, Montréal

<sup>†</sup> Canada CIFAR AI Chair

Many works have then sought to find structural assumptions and algorithmic modifications that accelerate this sublinear rate of  $\mathcal{O}(1/K)$ . The strong convexity of the set (or more generally uniform convexity, see [Kerdreux et al. 2020]) is one of such structural assumptions which lead to various accelerated convergence rates, like linear convergence rates when the unconstrained optimum is outside the constraint set [Levitin et al. 1966; Demyanov et al. 1970; Dunn 1979; Rector-Brooks et al. 2019] or sublinear rates  $\mathcal{O}(1/K^2)$  when the function is also strongly convex but without restrictions on the position of the optimum [Garber et al. 2015]. However, to the best of our knowledge, there exists no affine invariant analysis for these accelerated regimes stemming from the strong convexity of the constraint set  $\mathcal{C}$ .

In these “non affine invariant” analyses, structural assumptions like the  $L$ -smoothness (Definition 1.2) of  $f$  and the  $\alpha$ -strong convexity of  $\mathcal{C}$  (Definition 1.3) lead to accelerated convergence rate of the Frank-Wolfe algorithm, but are typically conditioned on parameters  $L, \alpha$  and others, which depend on a particular choice of a norm. This is surprising given that the Frank-Wolfe algorithm (under appropriate line-search) does not depend on any norm choice.

Recall that the smoothness of a function and the strong convexity of a set are defined as follows.

**Definition 1.2** *The function  $f$  is **smooth** over the set  $\mathcal{C}$  w.r.t. the norm  $\|\cdot\|$  if there exists a constant  $L > 0$  such that, for any  $x, y \in \mathcal{C}$ , we have*

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|x - y\|^2. \quad (3)$$

**Definition 1.3** *A set  $\mathcal{C}$  is  **$\alpha$ -strongly convex** with respect to a norm  $\|\cdot\|$  if, for any  $(x, y) \in \mathcal{C}$ ,  $\gamma \in [0, 1]$  and  $\|z\| \leq 1$ , we have*

$$\gamma x + (1 - \gamma)y + \alpha\gamma(1 - \gamma)\|x - y\|^2 z \in \mathcal{C}. \quad (4)$$

Obtaining *practical* accelerated affine invariant rates is hard, as an affine invariant step size is required. Indeed, some adaptive step sizes rely on theoretical affine invariant quantities which are in general not accessible. Therefore, by practical, we consider rates that can be achieved without a deep knowledge of the problem structure and constants.

For instance, scheduled step sizes, e.g.  $\gamma_k = \frac{2}{k+2}$ , makes the Frank-Wolfe algorithm practically affine covariant, yet they do not capture the accelerated convergence regimes. Exact line-search guarantees a practically affine covariant algorithm while capturing accelerated convergence regimes but significantly increases the time to perform a single iteration. Finally, it is possible to use backtracking line-search

such as [Pedregosa et al. 2020]. Unfortunately, backtracking techniques rely on the choice of a specific norm, thus breaking affine invariance of the algorithm. This raises naturally the following questions:

*Can we derive affine invariant rates for the Frank-Wolfe algorithm on strongly convex sets?*

*Can we design an affine invariant backtracking line-search for Frank-Wolfe algorithms?*

This work provides a positive answer to these questions, by proposing the following contributions.

**Contributions.** In this paper, **1)** we conduct affine invariant analyses of the Frank-Wolfe Algorithm 1, when the function  $f$  is smooth w.r.t. to a specific distance function  $\omega(\cdot)$  and the set  $\mathcal{C}$  is strongly convex also w.r.t.  $\omega(\cdot)$ . We then introduce new structural assumptions extending the class of problems for which such accelerated regimes hold in the case of Frank-Wolfe, called *directionally smooth functions with direction  $\delta$* . From this definition, **2)** we propose an affine invariant backtracking line-search for finding the optimal step size. Finally, **3)** we show that existing backtracking line-search methods, which use a specific norm, converges surprisingly to the optimal norm-invariant, affine invariant step size, meaning that affine-dependent and affine invariant backtracking techniques perform similarly.

**Outline.** In Section 2, we motivate the need for affine invariant analysis of Frank-Wolfe on strongly convex sets. In Section 3 and 4, we introduce the structural assumptions on the optimization problem that we will consider for analysing Frank-Wolfe. In Section 5 we detail our affine invariant analysis of Frank-Wolfe on strongly convex set. In Section 6 and 7 we provide a backtracking line-search that directly estimate the affine invariant quantities we developed and we explain how it relates with existing ones. We conclude in Section 8 with numerical experiments.

**Related Work.** Other linear convergence rates of Frank-Wolfe algorithms exists with affine invariant analysis. For instance, corrective variants of Frank-Wolfe exhibit (affine invariant) linear convergence rates when the constraint set is a polytope [Lacoste-Julien et al. 2013; Lacoste-Julien et al. 2015a] and the objective function is (generally) strongly convex. See Table 1 for a review of all affine invariant analyses of Frank-Wolfe algorithms.

Related Work	$\mathcal{C}$	Str. cvx. $f$	$x^*$	Algo	Step size	Rate
Clarkson [2010]	Simplex	$\times$	Any	FW	Scheduled	$\mathcal{O}(1/K)$
Jaggi [2013]	Convex	$\times$	Any	FW	Scheduled	$\mathcal{O}(1/K)$
Lacoste-Julien et al. [2013]	Any	$\checkmark$	Interior	FW	Exact ls	Linear
Lacoste-Julien et al. [2015a]	Polytope	$\checkmark$	Any	Corr. FW	Exact ls	Linear
Gutman et al. [2020]						
<b>Our work</b>	Strongly cvx	$\times$	$\nabla f(x^*) \neq 0$	FW	Backtracking ls	Linear
	Strongly cvx	$\checkmark$	Any	FW	Backtracking ls	$\mathcal{O}(1/K^2)$

Table 1: Existing *affine invariant* analysis of Frank-Wolfe for smooth convex functions under different schemes. **Strong convexity.** The strong convexity assumption is to be taken in a broad sense. In [Lacoste-Julien et al. 2013; Lacoste-Julien et al. 2015a], the authors consider “generalized geometric strong convexity” (see their Eq. 39), an affine invariant measure of (generalized) strong convexity, while [Gutman et al. 2020] consider strongly convex functions relative to a pair  $(\mathcal{C}, \omega)$  where  $\omega$  is a distance-like function. In our work, we do not directly assume strong convexity, but the *directional smoothness* of the function (see later Definition 4.1), whose constant is bounded if various assumptions are satisfied for problem (1) (Theorem 4.4). **Step size.** By *scheduled* step sizes, we consider, for instance, the classical  $\gamma_k = \frac{2}{k+2}$ . We denote by *exact-line search* when the optimal step size depends on an unknown affine invariant quantity, whose accessible upper-bounds are affine-dependent (thus breaking the affine invariance of FW).

These affine invariant analyses emphasize that there is no specific choice of norm to be made in Frank-Wolfe algorithms as well as there is no need for affine preconditioners. Frank-Wolfe algorithms are arguably *free-of-choice* methods, *i.e.* little needs to be known on the optimization problem’s structures to obtain the accelerated regimes. This is in line with recent works showing that the Frank-Wolfe methods exhibit accelerated adaptive behavior under a variety of structural constraints of (1) which depend on inaccessible parameters, *e.g.* Hölderian Error Bounds on  $f$  [Kerdreux et al. 2018a; Xu et al. 2018; Rinaldi et al. 2020] or local uniform convexity of  $\mathcal{C}$  [Kerdreux et al. 2020].

Affine invariant analyses introduce constants seeking to characterize structural properties without a specific choice of norm. This has then been the basis for works extending the accelerated convergence analysis to non-smooth or non-strongly convex functions [Pena 2019; Gutman et al. 2020], which then explore new structural assumptions on  $f$ .

## 2 “Affine-dependent” Analysis of FW

It is known that when the function is *smooth* (Definition 1.2), the set is *strongly-convex* (Definition 1.3) and the gradient is lower bounded  $\|\nabla f(x)\| \geq c$  over the constraint set (*i.e.*, the constraints are active), the Frank-Wolfe algorithm 1 converges linearly [Levitin et al. 1966; Demyanov et al. 1970; Dunn 1979], at rate (with  $h_k \triangleq f(x_k) - f_*$ )

$$h_k \leq \left( \max \left\{ \frac{1}{2}, 1 - \frac{c\alpha}{2L} \right\} \right)^k h_0. \quad (5)$$

Note that assuming the gradient to be lower bounded means the constraints are tight, *i.e.*, the solution of the unconstrained counterpart lies outside the set of constraints. However, the constants  $L$ ,  $\alpha$ , and  $c$  depend on the choice of the norm for the smoothness and the strong convexity. In contrast, the Frank-Wolfe algorithm and iterates do not depend on such a choice, due to its affine covariance. Therefore, the rate of Algorithm 1 should be affine invariant. Unfortunately, it is possible to show that the known theoretical analyses can be *arbitrarily* bad in the case where the constants  $L$ ,  $c$ ,  $\alpha$  depend on “affine variant” norms.

**Example 2.1** Consider the projection problem

$$\min_x f(x) \stackrel{\text{def}}{=} \frac{1}{2} \|x - \bar{x}\|^2 \quad \text{such that } \frac{1}{2} \|x\|^2 \leq 1.$$

In such case, we have that  $L = 1$ ,  $\alpha = \frac{1}{\sqrt{2}}$  and  $c = 1 - \|\bar{x}\|$  ( $L$ ,  $\alpha$  and  $c$  are defined according to the  $\ell_2$  norm). However, if we transform the problem into  $\min_y f(By)$ , the new constants become

$$L = \sigma_{\max}(B), \quad \alpha = \frac{\sigma_{\min}(B)}{\sqrt{2}\sigma_{\max}(B)}, \quad c = \sigma_{\max}(B)(1 - \|\bar{x}\|).$$

Comparing the rate (5) of the two problems, identical to the eyes of the FW algorithm, we have that

$$\begin{aligned} f(x_k) - f^* &\leq \left( 1 - \frac{(1 - \|\bar{x}\|)}{2\sqrt{2}} \right)^k (f(x_0) - f^*), \\ f(By_k) - f^* &\leq \left( 1 - \frac{(1 - \|\bar{x}\|)}{2\sqrt{2}} \kappa^{-1}(B) \right)^k (f(x_0) - f^*), \end{aligned}$$

where  $\kappa(B) = \frac{\sigma_{\max}(B)}{\sigma_{\min}(B)}$  is the condition number of  $B$ . This means we can artificially make a large theoretical upper bound on the rate of convergence by using an

ill-conditioned transformation (i.e.,  $\kappa(B)$  large). However, the speed of convergence of FW iterates are not affected by any linear transformation (due to their affine-covariance), therefore the upper bound will not be representative of the true rate of convergence of FW.

When the optimum is in the relative interior of any compact set  $\mathcal{C}$ , FW converges linearly when  $f$  is strongly convex [Guélat et al. 1986; Lacoste-Julien et al. 2013]. On the other hand, linear convergence on strongly convex sets does not require strong convexity of  $f$  when the solution of the unconstrained problem lies outside the set [Demyanov et al. 1970]. Our paper hence focuses on extending the analysis where the unconstrained optimum is outside the constrain set [Demyanov et al. 1970].

These two analysis cover most practical cases, but not the situation where the unconstrained optimum is close to the boundary of  $\mathcal{C}$ . A recent analysis on strongly convex sets of [Garber et al. 2015] is not restrictive w.r.t. the position of the unconstrained optimum but conservative (convergence rate of  $\mathcal{O}(1/K^2)$ ). It is interesting as it not only deals with the (previously unknown) situation where the unconstrained optimum is on the boundary on  $\mathcal{C}$ , but also when it is arbitrarily close to it, leading to poorly conditioned linear convergence regimes. In Appendix D, we provide an affine invariant analysis of [Garber et al. 2015].

### 3 Smoothness and Strong Convexity w.r.t. General Distance Functions

The major limitation in the definition of smoothness of a function (Definition 1.2) and the strong convexity of a set (Definition 1.3) is the presence of the norm in their definition, whose constants may be dependent on affine transformation of the space (see Example 2.1). Technically, the notion of norm in the definition of smoothness and strong convexity of a function can be extended to the concept of distance-generating function, for instance using Bregman divergence [Bauschke et al. 2017; Lu et al. 2018] or gauge functions [d’Aspremont et al. 2018].

Although is it classical to use different distance-generating functions  $\omega$  (that satisfies Assumption 3.1 below) to characterize the smoothness of a function, we are not aware of such analysis for strongly convex sets. We believe that such analysis may exist, but for completeness we propose here an extension of the strong convexity of a set w.r.t. a distance function  $\omega$ .

**Assumption 3.1** *The function  $\omega(\cdot)$  satisfies*

- $\omega(x) = 0 \Leftrightarrow x = 0$ ,
- **Positivity:**  $\omega(x) \geq 0$ ,

- **Triangular Inequality:**  $\omega(x + y) \leq \omega(x) + \omega(y)$
- **Positive homogeneity:**  $\omega(\gamma x) = \gamma \omega(x)$ ,  $\gamma \geq 0$ ,
- **Bounded asymmetry:**  $\max_x \frac{\omega(x)}{\omega(-x)} \leq \kappa_\omega$ .

Since  $\omega(x)$  is convex by the triangle inequality, we define the dual distance

$$\omega_*(v) = \max_{x: \omega(x) \leq 1} \langle v, x \rangle. \quad (6)$$

**Remark 3.2** *Usually, extensions of smoothness of a function use Bregman divergences (see e.g. [Lu et al. 2018; Bauschke et al. 2017]). However, the assumption that the distance-generating function is positively homogeneous is crucial in our analysis, which is unfortunately, not satisfied for most Bregman divergences.*

A typical example satisfying such assumptions are gauge functions, also called *Minkowski functional*,

$$\omega_{\mathcal{Q}}(v) \stackrel{\text{def}}{=} \underset{\tau \geq 0}{\operatorname{argmin}} \tau \quad \text{subject to } v \in \tau \mathcal{Q},$$

where  $0 \in \operatorname{int} \mathcal{Q}$ . Such distance-generating function satisfies Assumption 3.1 if the set  $\mathcal{Q}$  is convex and compact, and contains 0 in its interior. Moreover, gauge functions are affine invariant.

Usually, most works using gauge function assume that the set  $\mathcal{Q}$  is *centrally symmetric* [d’Aspremont et al. 2018; Molinaro 2020], which add the assumption that

$$\omega(x) = \omega(-x).$$

In that case, the gauge function is a norm [Rockafellar 1970, Theorem 15.2.]. Removing symmetry extends non-trivially the definition of strongly convex sets w.r.t. the distance function  $\omega$ . We now recall the definitions of smoothness and strong convexity of a function w.r.t. a distance function  $\omega$ .

**Definition 3.3** *A function  $f$  is smooth (resp. strongly convex) w.r.t. the distance function  $\omega$  if, for a constant  $L_\omega$  (resp.  $\mu_\omega$ ), the function satisfies*

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L_\omega}{2} \omega^2(y - x), \quad (7)$$

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu_\omega}{2} \omega^2(y - x). \quad (8)$$

**Definition 3.4** *A set  $\mathcal{C}$  is  $\alpha_\omega$ -strongly convex w.r.t.  $\omega$  if, for any  $(x, y) \in \mathcal{C}$  and  $\gamma \in [0, 1]$ , we have*

$$z_\gamma + \alpha_\omega \gamma (1 - \gamma) \frac{(1 - \gamma) \omega^2(x - y) + \gamma \omega^2(y - x)}{2} z \in \mathcal{C},$$

where  $z_\gamma = \gamma x + (1 - \gamma)y$ , for all  $z$  such that  $\omega(z) \leq 1$ .

This definition extends the one of strongly convex sets with a general distance function that may not be a norm, see for instance [Garber et al. 2015].

With Definition 3.4, the level sets of smooth and strongly convex functions are also strongly convex sets when the function  $\omega$  is used. Such results appear for instance in [Journée et al. 2010] when  $\omega$  is the  $\ell_2$  norm.

**Lemma 3.5 (Strong Convexity of Sets)** *Let  $f$  be a  $L$ -smooth and  $\mu$ -strongly convex function w.r.t.  $\omega$ . Then, the set*

$$\mathcal{C} = \{x : f(x) - f_* \leq R\}$$

*is  $\alpha$ -strongly convex w.r.t.  $\omega$ , with  $\alpha = \frac{\mu_\omega}{\kappa_\omega \sqrt{2L\omega R}}$ .*

We defer the proof in Appendix A. This result corresponds *exactly* to the one of [Journée et al. 2010, Theorem 12], when we use  $\omega = \|\cdot\|_2$ .

**Scaling Inequality.** All proofs of Frank-Wolfe methods on strongly convex sets leverage the same property. The *scaling inequality* (equivalent to strong convexity of  $\mathcal{C}$  [Goncharov et al. 2017, Theorem 2.1.]) crucially relates the Frank-Wolfe gap with  $\|x_t - v_t\|^2$ , see e.g. [Kerdreux et al. 2020, Lemma 2.1.]. We extend the scaling inequality to strongly convex sets with generic distance functions.

**Lemma 3.6 (Distance Scaling Inequality)**

*Assume  $\mathcal{C}$  is  $\alpha_\omega$ -strongly convex w.r.t.  $\omega$ . Then for any  $x \in \mathcal{C}$ ,  $\phi \in \mathbb{R}^d \setminus \{0\}$ , and  $v_\phi \in \operatorname{argmax}_{v \in \mathcal{C}} \langle \phi, v \rangle$ , we have  $\phi \in N_{\mathcal{C}}(v_\phi)$  (normal cone) and*

$$\langle \phi, v_\phi - x \rangle \geq \frac{\alpha_\omega}{2} \omega_*(\phi) \omega^2(v_\phi - x). \quad (9)$$

*In particular for any iterate  $x_k$  of Frank-Wolfe and its Frank-Wolfe vertex  $v_k$  (Line 2 in Algorithm 1), we have*

$$\langle -\nabla f(x_k); v_k - x_k \rangle \geq \frac{\alpha_\omega}{2} \omega_*(-\nabla f(x_k)) \omega^2(v_k - x_k).$$

**Proof.** We start with  $v_\phi = \operatorname{argmax}_{v \in \mathcal{C}} \langle \phi; v \rangle$ . Then, we use the definition of strong convexity of a set,

$$\gamma x + (1-\gamma)v_\phi + \alpha_\omega \gamma(1-\gamma)D_\gamma(x-v_\phi)z \in \mathcal{C} \quad \forall z : \omega(z) \leq 1.$$

where  $D_\gamma(x-y) \stackrel{\text{def}}{=} \frac{\gamma \omega^2(x-y) + (1-\gamma)\omega^2(y-x)}{2}$ . Then, by optimality of  $v_\phi$ ,

$$\langle \phi; v_\phi \rangle \geq \langle \phi; \gamma x + (1-\gamma)v_\phi + \alpha_\omega \gamma(1-\gamma)D_\gamma(x-v_\phi)z \rangle$$

After simplification,

$$\langle \phi; v_\phi - x \rangle \geq \alpha_\omega(1-\gamma)D_\gamma(x-v_\phi)\langle \phi; z \rangle$$

which holds in particular when  $\phi = -\nabla f(x)$ ,  $\gamma = 0$  and  $z$  being the argmax (see (6)). ■

## 4 Directional Smoothness

We separately introduced smoothness for functions, and strong convexity for sets w.r.t. a distance function  $\omega$ . Analyses of Frank-Wolfe algorithm on strongly convex sets [Levitin et al. 1966; Demyanov et al. 1970; Dunn 1979] show that, when  $f$  is convex and smooth, and the unconstrained minima of  $f$  are outside of  $\mathcal{C}$ , there is linear convergence.

We hence propose a novel condition that mingles the smoothness of  $f$  with the strong convexity of  $\mathcal{C}$  when moving in a specific direction  $\delta$ . We are interested in particular with the FW direction and we will see later that this assumption guarantees a linear convergence rate in this case. We call this condition the *directional smoothness*.

**Definition 4.1** *The function  $f$  is directionally smooth with direction function  $\delta : \mathcal{C} \rightarrow \mathbb{R}^d$  if there exists a constant  $\mathcal{L}_{f,\delta} > 0$  s.t.  $\forall x \in \mathcal{C}$  and  $h > 0$  with  $x + h\delta(x) \in \mathcal{C}$ ,*

$$f(x + h\delta(x)) \leq f(x) - h\langle -\nabla f(x), \delta(x) \rangle + \frac{\mathcal{L}_{f,\delta} h^2}{2} \langle -\nabla f(x), \delta(x) \rangle. \quad (10)$$

The rationale of Definition 10 is to replace the norm in the usual smoothness condition (Definition 1.2) by a scalar product between the *direction* and the negative gradient, in order to get an affine invariant quantity for the FW direction (see Proposition 4.3 below).

Assuming  $\delta(x)$  is a descent direction, i.e.,  $\langle -\nabla f(x), \delta(x) \rangle > 0$ , we can obtain a minimization algorithm for  $f$ , by minimizing (10) over  $h$ ,

$$x_{k+1} = x_k + h_{\text{opt}} \delta(x_k), \quad h_{\text{opt}} = \min\{h_{\text{max}}; \mathcal{L}_{f,\delta}^{-1}\}.$$

**Example 4.2** *(Gradient descent on smooth functions) The gradient algorithm uses  $\delta(x) = -\nabla f(x)$ . In such case, the function is directionally smooth with constant  $L$ , and we obtain*

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) - h\|\nabla f(x)\|^2 + \frac{Lh^2}{2}\|\nabla f(x)\|^2 \\ &= f(x) + h\left(\frac{Lh}{2} - 1\right)\|\nabla f(x)\|^2. \end{aligned}$$

*The best  $h$  is given by  $h_{\text{opt}} = \frac{1}{L}$ , which is also the optimal one [Nesterov 2013].*

The advantage of directional smoothness is its affine invariance in the case where  $\delta(x)$  is the FW step.

**Proposition 4.3 (Affine Invariance of  $\mathcal{L}_{f,\delta}$ )**

*If  $\delta(x)$  is affine covariant (e.g. the FW direction  $\delta(x) \triangleq v(x) - x$ ), then  $\mathcal{L}_{f,\delta}$  in (10) is invariant to an affine transformation of the constraint set (proof in Appendix B.2).*

The next theorem shows that, in the case of the FW algorithm, the directional smoothness constant is bounded if the function is smooth and the set is strongly convex for any distance function  $\omega$ . We use this result later, to show that affine invariant backtracking line-search is equivalent to using the best distance function  $\omega$  to define  $L_\omega$ ,  $c_\omega$  and  $\alpha_\omega$ .

**Theorem 4.4 (Directional Smoothness of FW)**

Consider the function  $f$ , smooth w.r.t. the distance function  $\omega$ , with constant  $L_\omega$ , and the set  $\mathcal{C}$ , strongly convex with constant  $\alpha_\omega$ .

Let  $\delta(x) = x - v(x)$ ,  $v(x)$  being the FW corner

$$v(x) \stackrel{\text{def}}{=} \operatorname{argmin}_{v \in \mathcal{C}} \langle \nabla f(x), v \rangle.$$

Then, if  $\omega_*(-\nabla f(x)) > c_\omega$  for all  $x \in \mathcal{C}$  and some  $c_\omega > 0$ , the function  $f(x)$  is directionally smooth w.r.t. to  $\omega$ , with constant

$$\mathcal{L}_{f,\delta} \leq \frac{L_\omega}{c_\omega \alpha_\omega}. \quad (11)$$

**Proof.** See Appendix A.1 for the proof. ■

## 5 Affine Invariant Linear Rates

With the directional smoothness constant  $\mathcal{L}_{f,\delta}$  (affine invariant when  $\delta$  is the FW direction), Theorem 5.1 shows an affine invariant linear rate of convergence of FW, generalizing existing convergence results of Frank-Wolfe on strongly convex sets [Levitin et al. 1966; Demyanov et al. 1970; Dunn 1979].

**Theorem 5.1 (Affine Invariant Linear Rates)**

Assume  $f$  is a convex function and directionally smooth with direction function  $\delta$  with constant  $\mathcal{L}_{f,\delta}$ . Then, the FW Algorithm 1 with step size

$$h_{\text{opt}} = \min \left\{ 1, \frac{1}{\mathcal{L}_{f,\delta}} \right\}, \quad \text{with } \delta = v(x) - x,$$

or with line-search, where  $v(x)$  is the FW corner

$$v(x) = \operatorname{argmin}_{v \in \mathcal{C}} \langle \nabla f(x), v \rangle,$$

converges linearly, at rate

$$f(x_k) - f_\star \leq \max \left\{ \frac{1}{2}, 1 - \frac{1}{2\mathcal{L}_{f,\delta}} \right\} (f(x_{k-1}) - f_\star).$$

**Proof.** We start with the directional smoothness assumption. For  $0 < h < 1$ ,

$$f(x_{k+1}) \leq f(x_k) + \left( h - \frac{\mathcal{L}_{f,\delta} h^2}{2} \right) \langle \nabla f(x_k), \delta(x_k) \rangle$$

After minimization, we have two possibilities:  $h_{\text{opt}} = \frac{1}{\mathcal{L}_{f,\delta}}$  or  $h_{\text{opt}} = 1$ . In the first case, we obtain

$$f(x_{k+1}) \leq f(x_k) + \frac{1}{2\mathcal{L}_{f,\delta}} \langle \nabla f(x_k), \delta(x_k) \rangle$$

Notice that the scalar product in the right-hand-side is the negative dual gap of Frank-Wolfe, that satisfies

$$\langle \nabla f(x_k), v(x) - x \rangle \leq -(f(x_k) - f_\star),$$

which gives the desired result. The second case follows immediately. ■

This provides an affine invariant analysis of the linear convergence regimes of FW on strongly convex sets.

The next proposition shows that the directional constant in Theorem 5.1 is bounded by (11) w.r.t. the distance function  $\omega$  that gives the best ratio. This means that the Frank-Wolfe method acts like it optimizes the function in the best possible geometry, i.e., the geometry that gives the *best constants*.

**Proposition 5.2 (Optimality of Dir. Smoothness)**

Let  $\Omega$  the set of function defined as

$$\Omega = \{ \omega : \omega \text{ satisfies assumptions 3.1} \}.$$

Then, the directional smoothness constant follows

$$\mathcal{L}_{f,\delta} \leq \min_{\omega \in \Omega} \frac{L_\omega}{c_\omega \alpha_\omega},$$

where  $L_\omega$  is the smoothness constant of the function  $f$ ,  $\alpha_\omega$  the strong convexity of the set  $\mathcal{C}$  and

$$c_\omega \leq \omega_*(-\nabla f(x)), \quad \forall x \in \mathcal{C}.$$

**Proof.** The proof is immediate by noticing that the FW algorithm do not use  $\omega$ , therefore we can choose the best  $\omega$  in Theorem 4.4. ■

To obtain a similar affine invariant analysis without restriction on the position of the optimum, i.e. the  $\mathcal{O}(1/K^2)$  analysis in [Garber et al. 2015], one can define a similar property to the direction smoothness defined in Section 4. This new structural assumption additionally mingles together with the strong convexity of  $f$ . We provide details in Appendix D. We choose to focus the analysis for the linear convergence in the main text as it is the one most significant in practice.

## 6 Affine Invariant Backtracking

In previous sections, we proposed new constants to bound the rate of convergence of the Frank-Wolfe algorithm, which is affine invariant. The significant advantage of these constants is that, like FW, they are

independent of any norm. However, the optimal step size of Frank-Wolfe needs the knowledge of these constants.

We propose in this section an affine invariant backtracking technique (Algorithm 2), based on directional smoothness. By construction, the backtracking technique finds automatically an estimate of the directional smoothness that satisfies

$$\mathcal{L}_k < 2\mathcal{L}_{f,\delta}, \quad k \geq \log_2 \left( \frac{\mathcal{L}_0}{\mathcal{L}_{f,\delta}} \right).$$

---

**Algorithm 2** Affine invariant backtracking

**Input:** FW corner  $v_k$ , point  $x_k$ , directional smoothness estimate  $\mathcal{L}_k$ , function  $f$ .

- 1:  $\mathcal{L} \leftarrow \mathcal{L}_k$ . Define the optimal step size and next iterate in the function of the directional Lipschitz constant:

$$\begin{aligned} \gamma_*(\mathcal{L}) &\stackrel{\text{def}}{=} \min\left\{\frac{1}{\mathcal{L}}, 1\right\}, \\ x(\mathcal{L}) &\stackrel{\text{def}}{=} (1 - \gamma_*(\mathcal{L}))x_k + \gamma_*(\mathcal{L})v_k. \end{aligned}$$

- 2: Create the model of  $f$  between  $x_k$  and  $x(\mathcal{L})$  based on equation (10),

$$m(\mathcal{L}) \stackrel{\text{def}}{=} f(x_k) + \gamma_*(\mathcal{L}) (1 - \gamma_*(\mathcal{L})) \langle \nabla f(x_k), v_k - x_k \rangle$$

- 3: Set the current estimate  $\tilde{\mathcal{L}} \stackrel{\text{def}}{=} \frac{\mathcal{L}_k}{2}$ .
- 4: **while**  $f(x(\tilde{\mathcal{L}})) > m(\tilde{\mathcal{L}})$  (Sufficient decrease not met because  $\tilde{\mathcal{L}}$  is too small) **do**
- 5: Double the estimate :  $\tilde{\mathcal{L}} \leftarrow 2 \cdot \tilde{\mathcal{L}}$ .
- 6: **end while**

**Output:** Estimate  $\mathcal{L}_{k+1} = \tilde{\mathcal{L}}$ , iterate  $x_{k+1} = x(\tilde{\mathcal{L}})$

---

## 7 Why Backtracking FW with norms is so efficient?

The step size strategy in Frank-Wolfe usually drives its practical efficiency. Sometimes, setting the step size optimally w.r.t. the theoretical analysis may be suboptimal in practice. Recently, Pedregosa et al. [2020] analyze the rate of the Frank-Wolfe algorithm for smooth function, using *backtracking line search*, described in Algorithm 3, Appendix C.

Algorithm 3 in Appendix C is adaptive to the local smoothness constant, and ensures  $L_{k+1} < 2L_f$ ,  $L_f$  being the smoothness constant of the function in the  $\ell_2$  norm. Pedregosa et al. [2020] observed that the estimate of the Lipschitz constant is often significantly smaller than the theoretical one; they wrote: “We compared the average Lipschitz estimate  $L_t$  and the

$L$ , the gradient’s Lipschitz constant. We found that across all datasets the former was more than an order of magnitude smaller, highlighting the need to use a local estimate of the Lipschitz constant to use a large step size.”

With our analysis, however, we can explain why the estimate of the smoothness constant is much better than the theoretical one. The answer is simple:

*Despite using a non-affine invariant bound, the step size resulting from the estimation of the Lipschitz constant via the backtracking line-search finds  $\frac{1}{\mathcal{L}_{f,\delta}}$ .*

**Proposition 7.1** Consider the “local Lipschitz constant”  $L_{loc}(x)$  that satisfies (3) with  $y = x + h\delta(x)$ , i.e.,

$$\begin{aligned} f(x + h\delta(x)) &\leq f(x) + \nabla f(x)(x + h\delta(x)) \\ &\quad + L_{loc}(x) \frac{h^2}{2} \|\delta(x)\|_2^2. \end{aligned}$$

Then,  $L_{loc}(x)$  is bounded by

$$L_{loc}(x) \leq \mathcal{L}_{f,\delta} \frac{\langle -\nabla f(x), \delta(x) \rangle}{\|\delta(x)\|_2^2}.$$

Assuming  $L_{loc}(x)$  “locally constant”, the backtracking line-search finds  $L_k < 2L_{loc}(x_k)$ , and its step size  $\gamma_*$  satisfies

$$\min \left\{ 1, \frac{1}{2\mathcal{L}_{f,\delta}} \right\} \leq \gamma_*.$$

**Proof.** See Appendix B.1 for the proof. ■

Therefore, the optimal step size from the backtracking line-search with the  $\ell_2$  norm is *exactly* the optimal affine invariant step size of our affine invariant analysis from Theorem 5.1.

In conclusion, *even if we use non-affine invariant norms to find the smoothness constant, surprisingly, the backtracking procedure finds the optimal, affine invariant step size.*

## 8 Illustrative Experiments

**Quadratic / logistic regression.** We consider the constrained quadratic and logistic regression problem,

$$\min_{x \in \mathcal{C}} \frac{1}{n} \sum_{i=1}^n l(a_i^T x, y_i), \quad (12)$$

where  $l$  is the quadratic or the logistic loss. Here we adopt the  $\ell_2$ -ball, defined as

$$\mathcal{C} = \{x : \|x\|_2 \leq R\}, \quad R > 0.$$

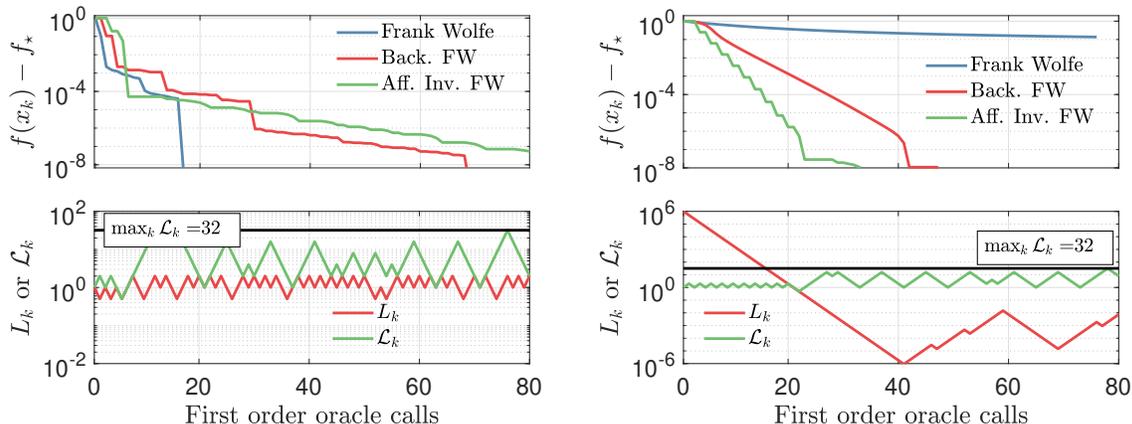


Figure 1: Comparison of FW variants on the projection problem. Left:  $B = I$ , Right:  $\kappa(B) = 10^6$ . The top row is the gap  $f_k - f^*$ , and the bottom row corresponds to the estimation of the directional-smoothness constant  $\mathcal{L}_k$  or the smoothness constant  $L_k$ , where the black line report the maximum value of  $\mathcal{L}_k$ . The reason why adaptive FW methods are slower in the left figure is because, in the worst case, the number of iterations to reach a certain precision can be up to four times larger than the worst-case bound on non-adaptive methods. We clearly see that the directional smoothness parameter  $\mathcal{L}_{f,\delta}$  is affine invariant, as its estimate is  $\max_k \mathcal{L}_k = 32$  in both scenarios.

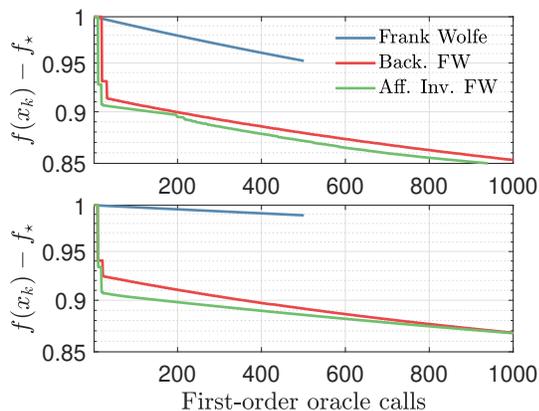


Figure 2: Classification problem on Madelon dataset, with (Top) Quadratic loss and (Bottom) Logistic loss.

Specifically, we compare our affine invariant backtracking method in Algorithm 2 against the naive FW Algorithm 1 with step size  $1/L$  [Demyanov et al. 1970] and back-tracking FW [Pedregosa et al. 2020] on the Madelon dataset [Guyon et al. 2007]. The results are shown in Figure 2. In detail, we set  $R$  such that the unconstrained optimum  $\mathbf{x}^*$  satisfies  $\|\mathbf{x}^*\|_2 = 1.1R$ , and the initial iterate  $\mathbf{x}_0 = \mathbf{0}$ . As predicted by our theory, the affine invariant algorithm performs well at the beginning, but after a few iterations the two backtracking techniques behave similarly.

**Projection.** We solve here the projection problem described in Example 2.1, for two cases of  $B$ : One that corresponds to the original problem, i.e.  $B = I$ , the

second one where  $B$  is an ill-conditioned matrix (with the condition number  $\kappa(B) = 10^6$ ). The vector  $x_0$  is random in the  $\ell_2$  ball, and  $\bar{x} = \mathbf{1}_d \cdot (1.1/\sqrt{d})$ . We report the results in Figure 1. We compare the standard FW algorithm with step size  $1/L$ , the FW with backtracking line-search (Algorithm 3) and FW with affine invariant backtracking technique (Algorithm 2). If the problem is well-conditioned ( $\kappa(B) = 1$ ), all methods perform similarly. This is not the case, however, for the ill-conditioned setting, where the FW with no adaptive step size converges extremely slowly compared to the two other methods. We also see that the affine invariant backtracking converges quicker than the standard backtracking. This is explained by the fact that the latter takes a longer time to find the right constant  $L_k$ , while  $\mathcal{L}_k$  remains untouched after an affine transformation.

## 9 Conclusion

In this paper, our theoretical convergence results on strongly convex sets complete the series of accelerated affine invariant analyses of Frank-Wolfe algorithms. To obtain these, we formulate a new structural assumption with respect to general distance functions, the directional smoothness, which we will explore more systematically in future works. Also, we present a new affine invariant backtracking line-search method based on directional smoothness. Within our framework of analysis, we provide a new explanation for the reasons behind the efficiency of the existing backtracking line search, and we show theoretically and experimentally they also find affine-invariant step sizes.

## Acknowledgments

This research was partially supported by the Canada CIFAR AI Chair Program. Simon Lacoste-Julien is a CIFAR Fellow in the Learning in Machines & Brains program.

## References

- Alayrac, Jean-Baptiste et al. (2016). “Unsupervised learning from narrated instruction videos”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4575–4583.
- Bauschke, Heinz H, Jérôme Bolte, and Marc Teboulle (2017). “A descent lemma beyond Lipschitz gradient continuity: first-order methods revisited and applications”. In: *Mathematics of Operations Research* 42.2, pp. 330–348.
- Bojanowski, Piotr et al. (2014). “Weakly supervised action labeling in videos under ordering constraints”. In: *European Conference on Computer Vision*. Springer, pp. 628–643.
- Braun, Gábor, Sebastian Pokutta, and Daniel Zink (2017). “Lazifying Conditional Gradient Algorithms”. In: *Proceedings of ICML*.
- Braun, Gábor et al. (2018). “Blended Conditional Gradients: the unconditioning of conditional gradients”. In: *arXiv preprint arXiv:1805.07311*.
- Carderera, Alejandro and Sebastian Pokutta (2020). “Second-order Conditional Gradients”. In: *arXiv preprint arXiv:2002.08907*.
- Clarkson, K.L. (2010). “Coresets, sparse greedy approximation, and the Frank-Wolfe algorithm”. In: *ACM Transactions on Algorithms (TALG)* 6.4, p. 63.
- Combettes, Cyrille W., Christoph Spiegel, and Sebastian Pokutta (2020a). “Projection-Free Adaptive Gradients for Large-Scale Optimization”. In: eprint: [arXiv:2009.14114](https://arxiv.org/abs/2009.14114).
- Combettes, Cyrille W and Sebastian Pokutta (2020b). “Boosting Frank-Wolfe by Chasing Gradients”. In: *arXiv preprint arXiv:2003.06369*.
- Courty, Nicolas et al. (2016). “Optimal transport for domain adaptation”. In: *IEEE transactions on pattern analysis and machine intelligence* 39.9, pp. 1853–1865.
- d’Aspremont, Alexandre, Cristobal Guzman, and Martin Jaggi (2018). “Optimal affine-invariant smooth minimization algorithms”. In: *SIAM Journal on Optimization* 28.3, pp. 2384–2405.
- Demyanov, V. F. and A. M. Rubinov (1970). “Approximate Methods in Optimization Problems”. In: *Modern Analytic and Computational Methods in Science and Mathematics*.
- Dunn, Joseph C (1979). “Rates of convergence for conditional gradient algorithms near singular and non-singular extremals”. In: *SIAM Journal on Control and Optimization* 17.2, pp. 187–211.
- Frank, Marguerite, Philip Wolfe, et al. (1956). “An algorithm for quadratic programming”. In: *Naval research logistics quarterly* 3.1-2, pp. 95–110.
- Garber, Dan and Elad Hazan (2015). “Faster rates for the frank-wolfe method over strongly-convex sets”. In: *32nd International Conference on Machine Learning, ICML 2015*.
- Goncharov, Vladimir V and Grigorii E Ivanov (2017). “Strong and weak convexity of closed sets in a Hilbert space”. In: *Operations research, engineering, and cyber security*. Springer, pp. 259–297.
- Guélat, Jacques and Patrice Marcotte (1986). “Some comments on Wolfe’s ‘away step’”. In: *Mathematical Programming*.
- Gutman, David H and Javier F Pena (2020). “The condition number of a function relative to a set”. In: *Mathematical Programming*, pp. 1–40.
- Guyon, Isabelle et al. (2007). “Competitive baseline methods set new standards for the NIPS 2003 feature selection benchmark”. In: *Pattern recognition letters* 28.12, pp. 1438–1444.
- Jaggi, Martin (2013). “Revisiting Frank-Wolfe: Projection-free sparse convex optimization”. In: *Proceedings of the 30th international conference on machine learning*. CONF, pp. 427–435.
- Journée, Michel et al. (2010). “Generalized power method for sparse principal component analysis.” In: *Journal of Machine Learning Research* 11.2.
- Kerdreux, Thomas, Alexandre d’Aspremont, and Sebastian Pokutta (2020). “Projection-Free Optimization on Uniformly Convex Sets”. In: eprint: [arXiv:2004.11053](https://arxiv.org/abs/2004.11053).
- Kerdreux, Thomas, Alexandre d’Aspremont, and Sebastian Pokutta (2018a). “Restarting Frank-Wolfe”. In: *arXiv preprint arXiv:1810.02429*.
- Kerdreux, Thomas, Fabian Pedregosa, and Alexandre d’Aspremont (2018b). “Frank-Wolfe with subsampling oracle”. In: *arXiv preprint arXiv:1803.07348*.
- Lacoste-Julien, Simon and Martin Jaggi (2013). “An affine invariant linear convergence analysis for Frank-Wolfe algorithms”. In: *arXiv preprint arXiv:1312.7864*.
- (2015a). “On the Global Linear Convergence of Frank-Wolfe Optimization Variants”. In: *Advances in Neural Information Processing Systems*. Ed. by C. Cortes et al. Vol. 28. Curran Associates, Inc., pp. 496–504. [arXiv:1511.05932v1](https://arxiv.org/abs/1511.05932v1) [[math.OC](https://arxiv.org/abs/1511.05932v1)].
- Lacoste-Julien, Simon, Fredrik Lindsten, and Francis Bach (2015b). “Sequential kernel herding: Frank-Wolfe optimization for particle filtering”. In: *arXiv preprint arXiv:1501.02056*.

- Levitin, Evgeny S and Boris T Polyak (1966). “Constrained minimization methods”. In: *USSR Computational mathematics and mathematical physics* 6.5, pp. 1–50.
- Lu, Haihao, Robert M Freund, and Yurii Nesterov (2018). “Relatively smooth convex optimization by first-order methods, and applications”. In: *SIAM Journal on Optimization* 28.1, pp. 333–354.
- Luise, Giulia et al. (2019). “Sinkhorn Barycenters with Free Support via Frank-Wolfe Algorithm”. In: *Advances in Neural Information Processing Systems*, pp. 9318–9329.
- Miech, Antoine, Ivan Laptev, and Josef Sivic (2018). “Learning a text-video embedding from incomplete and heterogeneous data”. In: *arXiv preprint arXiv:1804.02516*.
- Molinaro, Marco (2020). “Curvature of Feasible Sets in Offline and Online Optimization”. In: eprint: [arXiv:2002.03213](https://arxiv.org/abs/2002.03213).
- Mortagy, Hassan, Swati Gupta, and Sebastian Pokutta (2020). “Walking in the Shadow: A New Perspective on Descent Directions for Constrained Minimization”. In: eprint: [arXiv:2006.08426](https://arxiv.org/abs/2006.08426).
- Nesterov, Yurii (2013). *Introductory lectures on convex optimization: A basic course*. Vol. 87. Springer Science & Business Media.
- Paty, François-Pierre and Marco Cuturi (2019). “Subspace robust wasserstein distances”. In: *arXiv preprint arXiv:1901.08949*.
- Pedregosa, Fabian et al. (2020). “Linearly convergent Frank-Wolfe with backtracking line-search”. In: *International Conference on Artificial Intelligence and Statistics*. PMLR, pp. 1–10.
- Pena, Javier (2019). “Generalized conditional subgradient and generalized mirror descent: duality, convergence, and symmetry”. In: *arXiv preprint arXiv:1903.00459*.
- Peyre, Julia et al. (2017). “Weakly-supervised learning of visual relations”. In: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 5179–5188.
- Rector-Brooks, Jarrid, Jun-Kun Wang, and Barzan Mozafari (2019). “Revisiting projection-free optimization for strongly convex constraint sets”. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 33, pp. 1576–1583.
- Rinaldi, Francesco and Damiano Zeffiro (2020). “A unifying framework for the analysis of projection-free first-order methods under a sufficient slope condition”. In: eprint: [arXiv:2008.09781](https://arxiv.org/abs/2008.09781).
- Rockafellar, R Tyrrell (1970). *Convex analysis*. 28. Princeton university press.
- Seguin, Guillaume et al. (2016). “Instance-level video segmentation from object tracks”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3678–3687.
- Xu, Yi and Tianbao Yang (2018). “Frank-Wolfe Method is Automatically Adaptive to Error Bound Condition”. In: eprint: [arXiv:1810.04765](https://arxiv.org/abs/1810.04765).

## A Strong Convexity of Sets with asymmetric distance functions

Before presenting the proof, we introduce the following results, extending known properties from smooth and strongly convex sets.

**Proposition A.1** *If  $f$  is strongly convex w.r.t. the distance function  $\omega$ , then for  $\gamma \in [0, 1]$  we have*

$$f(\gamma x + (1 - \gamma)y) + \mu\gamma(1 - \gamma) \frac{\gamma\omega^2(x - y) + (1 - \gamma)\omega^2(y - x)}{2} \leq \gamma f(x) + (1 - \gamma)f(y)$$

**Proof.** Let  $z_\gamma = \gamma x + (1 - \gamma)y$ . We start with the definition,

$$\begin{aligned} f(z_\gamma) + \langle \nabla f(z_\gamma), x - z_\gamma \rangle + \frac{\mu}{2}\omega^2(x - z_\gamma) &\leq f(x) \\ f(z_\gamma) + \langle \nabla f(z_\gamma), y - z_\gamma \rangle + \frac{\mu}{2}\omega^2(y - z_\gamma) &\leq f(y) \end{aligned}$$

After multiplying by  $\gamma$  and  $1 - \gamma$  and adding the two inequalities, we have

$$f(z_\gamma) + \mu \frac{\gamma\omega^2(x - z_\gamma) + (1 - \gamma)\omega^2(y - z_\gamma)}{2} \leq \gamma f(x) + (1 - \gamma)f(y)$$

Since  $\omega^2(x - z_\gamma) = (1 - \gamma)^2\omega^2(y - x)$ , and  $\omega^2(y - z_\gamma) = \gamma^2\omega^2(x - y)$ , we obtain the desired result. ■

**Proposition A.2** *If  $f$  is convex and smooth w.r.t. the distance function  $\omega$ , then it holds that*

$$\frac{1}{2L}\omega_*^2(\nabla f(x) - \nabla f(y)) \leq f(y) - f(x) - \langle \nabla f(x), y - x \rangle$$

where  $\omega_*$  is the dual of the function  $\omega$ , written

$$\omega_*(v) \stackrel{\text{def}}{=} \max_{s: \omega(s) \leq 1} \langle v, s \rangle.$$

In particular, Proposition A.2 implies that, if  $f$  has a minimum  $x_*$ , then

$$\frac{1}{2L}\omega_*^2(-\nabla f(y)) \leq f(y) - f(x_*) \tag{13}$$

**Proof.** Let the function  $\phi(y) = f(y) - \langle \nabla f(x), y \rangle$ . This function is, by construction, smooth. Moreover,  $\min_y \phi(y)$  is attained when  $y = x$ . Since the function is smooth,

$$\min_y \phi(y) \leq \min_y \phi(z) + \langle \nabla \phi(z), y - z \rangle + \frac{L}{2}\omega^2(y - z)$$

Let  $\beta u = y - z$ , where  $\omega(u) = 1$  and  $\beta \geq 0$ . Then,

$$\min_y \phi(y) \leq \min_{\beta, u} \phi(z) + \beta \langle \nabla \phi(z), u \rangle + \frac{\beta^2 L}{2}$$

The minimum can be split into two minimization problems,

$$\min_y \phi(y) \leq \phi(z) + \min_{\beta \geq 0} \left( \frac{\beta^2 L}{2} - \beta \max_{u: \omega(u) \leq 1} \langle -\nabla \phi(z), u \rangle \right).$$

By definition of the dual of  $\omega$ ,

$$\min_y \phi(y) \leq \phi(z) + \min_{\beta \geq 0} \left( \frac{\beta^2 L}{2} - \beta \omega_*(-\nabla \phi(z)) \right).$$

Now, we can solve over  $\beta$ , which gives us

$$\min_y \phi(y) \leq \phi(z) - \frac{1}{2L} \omega_*^2(-\nabla\phi(z)).$$

Replacing the minimum by  $\phi(x)$ , and  $\phi$  by its expression, we get

$$f(x) - \langle \nabla f(x), x \rangle \leq f(z) - \langle \nabla f(x), z \rangle - \frac{1}{2L} \omega_*^2(\nabla f(x) - \nabla f(z)).$$

After reorganization, we get the desired result. ■

We can now show that level sets of a smooth and strong convex function are strongly convex sets, when they use the distance function  $\omega$ .

**Proof. (Proof of Lemma 3.5.)**

Consider the set

$$\mathcal{C} = \{x : f(x) - f_* \leq R\}$$

Let  $x, y \in \mathcal{C}$ . Let  $z_\gamma = \gamma x + (1 - \gamma)y$ , and consider the point  $z_\gamma + u$ . We have that

$$\begin{aligned} f(z_\gamma + u) - f_* &\leq f(z_\gamma) - f_* + \langle \nabla f(z_\gamma), u \rangle + \frac{L}{2} \omega^2(u), \\ &\leq f(z_\gamma) - f_* + \omega(-u) \max_{v: \omega(v) \leq 1} \langle -\nabla f(z_\gamma), v \rangle + \frac{L}{2} \omega^2(u), \\ &= f(z_\gamma) - f_* + \omega(-u) \omega_* (-\nabla f(z_\gamma)) + \frac{L}{2} \omega^2(u), \\ &\leq f(z_\gamma) - f_* + \kappa_\omega \omega(u) \sqrt{2L(f(z_\gamma) - f_*)} + \frac{L}{2} \omega^2(u). \end{aligned}$$

Therefore, to satisfy  $f(z_\gamma + u) - f_* \leq R$ , we need to ensure that

$$\underbrace{f(z_\gamma) - f_* - R}_{=\omega} + \underbrace{\kappa_\omega \sqrt{2L(f(z_\gamma) - f_*)} \omega(u)}_{=\beta} + \frac{L}{2} \omega^2(u) \leq 0$$

Solving the problem in  $\omega(u)$  gives

$$\omega(u) \leq \frac{-\beta + \sqrt{\beta^2 - 2L\omega}}{L}$$

We have that

$$\beta^2 - 2L\omega = 2L((f(z_\gamma) - f_*)(\kappa_\omega^2 - 1) + R)$$

Therefore,

$$\omega(u) \leq \sqrt{2} \frac{-\kappa_\omega \sqrt{(f(z_\gamma) - f_*)} + \sqrt{(f(z_\gamma) - f_*)(\kappa_\omega^2 - 1) + R}}{\sqrt{L}}$$

However, since the function is strongly convex,

$$f(z_\gamma) - f_* \leq \underbrace{\gamma f(x) + (1 - \gamma)f(y) - f_*}_{\leq R} - \mu\gamma(1 - \gamma) \frac{\gamma\omega^2(x - y) + (1 - \gamma)\omega^2(y - x)}{2}$$

Let  $D_\gamma = \gamma(1 - \gamma) \frac{\gamma\omega^2(x - y) + (1 - \gamma)\omega^2(y - x)}{2}$ . The inequality now reads

$$f(z_\gamma) - f_* \leq R - \mu D_\gamma. \tag{14}$$

Therefore, the condition on  $\omega$  becomes

$$\omega(u) \leq \sqrt{2} \frac{-\kappa_\omega \sqrt{R - \mu D_\gamma} + \sqrt{(R - \mu D_\gamma)(\kappa_\omega^2 - 1) + R}}{\sqrt{L}}$$

which gives

$$\omega(u) \leq \frac{\kappa_\omega \sqrt{2}}{\sqrt{L}} \left( -\sqrt{R - \mu D_\gamma} + \sqrt{R - \left(1 - \frac{1}{\kappa_\omega^2}\right) \mu D_\gamma} \right) \quad (15)$$

To simplify the expression in parenthesis, we multiply and divide by the conjugate of the square roots to get:

$$\begin{aligned} \left( -\sqrt{R - \mu D_\gamma} + \sqrt{R - \left(1 - \frac{1}{\kappa_\omega^2}\right) \mu D_\gamma} \right) &= \frac{R - \left(1 - \frac{1}{\kappa_\omega^2}\right) \mu D_\gamma - (R - \mu D_\gamma)}{\sqrt{R - \mu D_\gamma} + \sqrt{R - \left(1 - \frac{1}{\kappa_\omega^2}\right) \mu D_\gamma}} \\ &\geq \frac{1}{\kappa_\omega^2 2\sqrt{R}}. \end{aligned}$$

We can thus strengthen the condition (15) to:

$$\omega(u) \leq \frac{\mu D_\gamma}{\kappa_\omega \sqrt{2LR}}.$$

As the definition of a strongly convex set requires  $\omega(u) \leq \alpha_\omega D_\gamma$ , we conclude that the level set is strongly convex with at least the constant  $\alpha_\omega = \frac{\mu}{\kappa_\omega \sqrt{2LR}}$ . ■

#### A.1 Proof of Theorem 4.4

**Theorem A.3** Consider the function  $f$ , smooth w.r.t. the distance function  $\omega$ , with constant  $L_\omega$ , and the set  $\mathcal{C}$ , strongly convex with constant  $\alpha_\omega$ .

Let  $\delta(x) = x - v(x)$ ,  $v(x)$  being the FW corner

$$v(x) \stackrel{\text{def}}{=} \operatorname{argmin}_{v \in \mathcal{C}} \langle \nabla f(x), v \rangle.$$

Then, if  $\omega_*(-\nabla f(x)) > c_\omega$  for all  $x \in \mathcal{C}$ , the function  $f(x)$  is directionally smooth w.r.t. to  $\omega$ , with constant

$$\mathcal{L}_{f,\delta} \leq \frac{L_\omega}{c_\omega \alpha_\omega}. \quad (16)$$

**Proof.** We start by the definition of smooth functions between  $x$  and  $h\delta(x)$  for the distance function  $\omega$ . We have for all  $0 \leq h \leq 1$

$$f(x + h\delta(x)) \leq f(x) + h \langle \nabla f(x), \delta(x) \rangle + \frac{h^2 L_\omega}{2} \omega^2(\delta(x))$$

Using the scaling inequality in (9),

$$\langle -\nabla f(x), \delta(x) \rangle \geq \alpha_\omega \omega_*(-\nabla f(x)) \omega(\delta(x))^2.$$

We hence obtain

$$f(x + h\delta(x)) \leq f(x) + h \langle \nabla f(x), \delta(x) \rangle - \frac{h^2 L_\omega}{2} \frac{\langle \nabla f(x), \delta(x) \rangle}{\alpha_\omega \omega_*(-\nabla f(x))}.$$

Since  $\omega_*(-\nabla f(x)) > c_\omega$  for all  $x \in \mathcal{C}$ ,

$$f(x + h\delta(x)) \leq f(x) + h \langle \nabla f(x), \delta(x) \rangle - \frac{h^2}{2} \frac{L_\omega}{\alpha_\omega c_\omega} \langle \nabla f(x), \delta(x) \rangle.$$

which is the definition of directional smoothness. ■

## B Missing proofs

### B.1 Proof of Proposition 7.1

**Proposition B.1** We define the “local Lipschitz constant”  $L_{\text{loc}}(x)$ , which satisfies

$$L_{\text{loc}}(x) \stackrel{\text{def}}{=} \mathcal{L}_{f,\delta} \frac{\langle -\nabla f(x), \delta(x) \rangle}{\|\delta(x)\|^2}.$$

Then, assuming that the local Lipschitz constant is “locally constant”, the backtracking line-search finds  $L_k \leq 2L_{\text{loc}}(x_k)$ , and its step size  $\gamma_\star$  satisfies

$$\min \left\{ 1, \frac{1}{2\mathcal{L}_{f,\delta}} \right\} \leq \gamma_\star.$$

**Proof.** We start with the definition of directional smoothness,

$$f(x + h\delta(x)) \leq f(x) + h\langle \nabla f(x), \delta(x) \rangle + [\mathcal{L}_{f,\delta} \langle -\nabla f(x), \delta(x) \rangle] \frac{h^2}{2}.$$

Writing  $1 = \frac{\|\delta(x)\|_2^2}{\|\delta(x)\|_2^2}$ , the upper bound becomes

$$f(x) + h\langle \nabla f(x), \delta(x) \rangle + \left[ \frac{\mathcal{L}_{f,\delta} \langle -\nabla f(x), \delta(x) \rangle}{\|\delta(x)\|_2^2} \right] \frac{h^2 \|\delta(x)\|_2^2}{2}.$$

Defining

$$L_{\text{loc}}(x) \triangleq \frac{\mathcal{L}_{f,\delta} \langle -\nabla f(x), \delta(x) \rangle}{\|\delta(x)\|_2^2},$$

we obtain

$$f(x_k + h\delta(x_k)) \leq f(x_k) + h\langle \nabla f(x_k), \delta(x_k) \rangle + L_{\text{loc}}(x_k) \frac{h^2 \|\delta(x_k)\|_2^2}{2}.$$

If we assume that  $L_{\text{loc}}(x_k)$  is approximately constant, then Algorithm 3 finds  $L_k \leq 2L_{\text{loc}}(x_k)$ . Finally, using the definition of  $\gamma_\star$  in Algorithm 3, we have

$$\begin{aligned} \gamma_\star &= \min \left\{ \frac{-\nabla f(x_k) \langle v_k - x_k \rangle}{L_{\text{loc}}(x_k) \|v_k - x_k\|^2}, 1 \right\} \\ &\geq \min \left\{ \frac{1}{2\mathcal{L}_{f,\delta}}, 1 \right\}. \end{aligned}$$

■

### B.2 Proof of Proposition 4.3

**Proposition B.2 (Affine Invariance)** If  $\delta(x)$  is affine covariant (e.g. the Frank-Wolfe direction  $\delta(x) \triangleq v(x) - x$ ), then the constant  $\mathcal{L}_{f,\delta}$  in (10) is affine invariant. In other words, let

$$\tilde{f}(\cdot) \triangleq f(B\cdot), \quad \tilde{\delta}_{\tilde{c}}(\cdot) \triangleq \delta_{B \cdot c}(\cdot),$$

then  $\mathcal{L}_{\tilde{f}, \tilde{\delta}_{\tilde{c}}} = \mathcal{L}_{f,\delta}$ .

**Proof.** We start with the definition of directional smoothness, but with  $x \rightarrow By$ . The upper bound reads

$$f(By) + \left( h - \frac{\mathcal{L}_{f,\delta} h^2}{2} \right) \langle \nabla f(By), \delta(By) \rangle$$

Since we assumed  $\delta(By)$  affine covariant,

$$\delta(By) = B\tilde{\delta}_{\tilde{c}}(y).$$

Therefore,

$$f(By) + \left( h - \frac{\mathcal{L}_{f,\delta} h^2}{2} \right) \langle B^T \nabla f(By), \tilde{\delta}_{\tilde{\mathcal{C}}}(y) \rangle$$

Since  $\nabla \tilde{f}(y) = B^T \nabla f(By)$ , we have

$$\tilde{f}(\tilde{y} + h\tilde{\delta}_{\tilde{\mathcal{C}}}(y)) \leq \tilde{f}(y) + \left( h - \frac{\mathcal{L}_{f,\delta} h^2}{2} \right) \langle \nabla \tilde{f}(y), \tilde{\delta}_{\tilde{\mathcal{C}}}(y) \rangle$$

This means the function  $\tilde{f}$  is directionally smooth with constant  $\mathcal{L}_{f,\delta}$ , which proves the statement.  $\blacksquare$

## C Backtracking Line Search for Frank-Wolfe Steps

---

**Algorithm 3** Backtracking line-search for smooth functions [Pedregosa et al. 2020]

---

**Input:** FW corner  $v_k$ , point  $x_k$ , smoothness estimate  $L_k$ , function  $f$ .

1: Create the optimal step size and next iterate in the function of the Lipchitz estimate

$$\begin{aligned} \gamma_*(L) &\stackrel{\text{def}}{=} \min \left\{ \frac{-\nabla f(x_k)(v_k - x_k)}{L\|v_k - x_k\|^2}, 1 \right\}. \\ x(L) &\stackrel{\text{def}}{=} (1 - \gamma_*(L)) + \gamma_*(L)v_k \end{aligned}$$

2: Quadratic model of  $f$  between  $x_k$  and  $x(L)$ ,

$$m(L) \stackrel{\text{def}}{=} f(x_k) + \langle \nabla f(x_k), x(L) - x_k \rangle + \frac{L}{2} \|x(L) - x_k\|^2$$

3: Set the current estimate  $\tilde{L} \stackrel{\text{def}}{=} \frac{L_k}{2}$ .

4: **while**  $f(x(\tilde{L})) > m(\tilde{L})$  (Sufficient decrease not met because  $\tilde{L}$  is too small) **do**

5: Double the estimate :  $\tilde{L} \leftarrow 2 \cdot \tilde{L}$ .

6: **end while**

**Output:** Estimate  $L_{k+1} = \tilde{L}$ , iterate  $x_{k+1} = x(\tilde{L})$

---

## D Affine Invariant Analysis without Restriction on Optimum Location

In this section, we propose a modification of the directional smoothness defined in Section 4. This new assumption is the basis to obtain an affine invariant analysis of Frank-Wolfe on a strongly convex set without restriction on the position of the unconstrained optimum of  $f$ , as recently proposed in Garber et al. [2015].

**Outline.** In Theorem D.2, we prove a  $\mathcal{O}(1/K^2)$  sublinear convergence rate as in [Garber et al. 2015] when the function is *modified directionally smooth* (Definition D.1). In Theorem D.4, we prove that when  $\mathcal{C}$  is strongly convex, and  $f$  is smooth and strongly convex, then  $f$  is *modified directionally smooth* for the Frank-Wolfe direction with an affine invariant constant leading to better conditioned convergence rates than in [Garber et al. 2015]. Finally, in Proposition D.5, we show that the constant of modified directional smoothness is affine invariant.

We now define a modification of directional smoothness. It is a structural assumption on  $f$  constrained on  $\mathcal{C}$  designed at gathering the strong convexity of  $\mathcal{C}$ , the smoothness, and the strong convexity of  $f$  into a single quantity.

**Definition D.1 (Modified Directional Smoothness)** Let  $x_0 \in \mathcal{C}$ . The function  $f$  is called *modified directionally smooth* with direction function  $\delta : \mathcal{C} \rightarrow \mathbb{R}^N$  if there exists a constant  $\tilde{\mathcal{L}}_{f,\delta}(x_0) > 0$  such that  $\forall x \in \mathcal{C}$ ,

$$f(x + h\delta(x)) \leq f(x) + h\langle \nabla f(x), \delta(x) \rangle - \frac{\tilde{\mathcal{L}}_{f,\delta}(x_0) h^2}{2} \langle \nabla f(x), \delta(x) \rangle \sqrt{\frac{f(x_0) - f^*}{f(x) - f^*}}, \quad (17)$$

for  $0 < h < 1$ .

Note that the dependence of  $x_0$  in the definition of the modified directional smoothness is an artifact to obtain a dimensionless constant  $\tilde{\mathcal{L}}_{f,\delta}(x_0)$ .

As in Section 5, the modified directional smoothness constant  $\tilde{\mathcal{L}}_{f,\delta}$  is affine invariant in the case where  $\delta$  is the FW direction. We now derive an affine invariant accelerated sublinear rate of convergence of Frank-Wolfe providing an affine invariant analysis of [Garber et al. 2015].

**Theorem D.2 (Affine Invariant Accelerated Sublinear Rates)** *Let  $x_0 \in \mathcal{C}$  and assume  $f$  is a convex function and modified directionally smooth with direction function  $\delta$  and constant  $\tilde{\mathcal{L}}_{f,\delta}(x_0)$ . Then, the iterates  $x_k$  for the Frank-Wolfe Algorithm 1 with step size*

$$h_{\text{opt}} = \min \left\{ 1, \frac{1}{\tilde{\mathcal{L}}_{f,\delta}(x_0)} \sqrt{\frac{f(x_k) - f^*}{f(x_0) - f^*}} \right\}, \quad \text{with } \delta = v(x) - x,$$

or with exact line-search, where  $v(x)$  is the Frank-Wolfe corner

$$v(x) = \underset{v \in \mathcal{C}}{\operatorname{argmin}} \langle \nabla f(x), v \rangle,$$

satisfy

$$f(x_k) - f^* \leq \frac{4(f(x_0) - f^*) \max\{1, 18\tilde{\mathcal{L}}_{f,\delta}^2(x_0)\}}{(k+2)^2} \quad \text{for } k \geq 0.$$

**Proof.** The proof is similar to that of Theorem 5.1. We hence start with the modified directional smoothness assumption on  $f$ . For  $0 < h < 1$ ,

$$f(x_{k+1}) \leq f(x_k) + \left( h - \frac{\tilde{\mathcal{L}}_{f,\delta} h^2}{2} \sqrt{\frac{f(x_0) - f^*}{f(x_k) - f^*}} \right) \langle \nabla f(x_k), \delta(x_k) \rangle \quad (18)$$

After minimizing over  $h$ , we have two possibilities. The case with exact line-search follows immediately after these two cases. In the following, we use the notation  $h_k \stackrel{\text{def}}{=} f(x_k) - f^*$  for the primal suboptimality at  $x_k$ , and  $g_k \stackrel{\text{def}}{=} \langle -\nabla f(x_k), \delta(x_k) \rangle$  for the Frank-Wolfe gap at  $x_k$  (and note that  $g_k \geq h_k$  by convexity).

**Case 1:**  $h_{\text{opt}} = \frac{1}{\tilde{\mathcal{L}}_{f,\delta}(x_0)} \sqrt{\frac{f(x_k) - f^*}{f(x_0) - f^*}}$ . In such case, we obtain (subtract  $f^*$  on both sides of the inequality)

$$h_{k+1} \leq h_k - \frac{1}{2\tilde{\mathcal{L}}_{f,\delta}} \sqrt{\frac{h_k}{h_0}} g_k,$$

and since the Frank-Wolfe gap  $g_k$  upper bounds the primal suboptimality, we obtain

$$h_{k+1} \leq h_k \left[ 1 - \frac{1}{2\tilde{\mathcal{L}}_{f,\delta} \sqrt{h_0}} \sqrt{h_k} \right].$$

**Case 2:** With  $h_{\text{opt}} = 1$ , we have

$$h_{k+1} \leq h_k + \left( 1 - \frac{\mathcal{L}_{f,\delta}}{2} \sqrt{\frac{h_0}{h_k}} \right) g_k.$$

In that case, we have that  $\frac{1}{\tilde{\mathcal{L}}_{f,\delta}(x_0)} \sqrt{\frac{h_k}{h_0}} \geq 1$ . Hence we obtain

$$h_{k+1} \leq h_k - \frac{1}{2} g_k \leq \frac{1}{2} h_k$$

Finally, we have the following recursive relation on the sequence of primal suboptimality ( $h_k$ ):

$$\begin{aligned} h_{k+1} &\leq h_k \cdot \max \left\{ \frac{1}{2}, 1 - \frac{1}{2\tilde{\mathcal{L}}_{f,\delta} \sqrt{h_0}} \sqrt{h_k} \right\} \\ &= h_k \cdot \max \left\{ \frac{1}{2}, 1 - M \sqrt{h_k} \right\}, \end{aligned} \quad (19)$$

with  $M \stackrel{\text{def}}{=} \frac{1}{2\tilde{\mathcal{L}}_{f,\delta}(x_0)\sqrt{h_0}}$ . The inequality (19) is exactly the same recurrence that was analyzed by Garber et al. [2015] (see their Equation (7), with the same notation for  $M$ ), where they have shown a  $\mathcal{O}(1/K^2)$  convergence rate. The exact constant is obtained by following the very same proof as [Garber et al. 2015], *i.e.* proving by induction that there exists  $C$  such that  $h_k \leq C/(k+2)^2$ . The base case  $k = 0$  can be trivially obtained by letting  $C \geq 4h_0$ .<sup>1</sup> Their induction step was shown by requiring that  $C \geq \frac{18}{M^2}$ . Thus using  $C = \max\{4h_0, \frac{18}{M^2}\}$  (and re-arranging) proves the statement of our theorem. ■

The following lemma will be used in the proof of the bound on the modified directional smoothness.

**Lemma D.3** *Consider a compact convex set  $\mathcal{C}$ . Assume  $f$  is a  $\mu_\omega$ -strongly convex function with respect to  $\omega$ . Let  $x^*$  be the minimum of  $f$  on  $\mathcal{C}$ . Then, for any  $x \in \mathcal{C}$ , we have*

$$\omega_*(\nabla f(x)) \geq \sqrt{\frac{\mu_\omega}{2}} \sqrt{f(x) - f(x^*)}. \quad (20)$$

**Proof.** Let  $x \in \mathcal{C}$ . From Definition 3.3, we have that

$$f(x) \geq f(x^*) + \langle \nabla f(x^*), x - x^* \rangle + \frac{\mu_\omega}{2} \omega^2(x - x^*).$$

Hence with the optimality conditions, *i.e.*  $\langle \nabla f(x^*), x - x^* \rangle \geq 0$ , we have

$$f(x) - f(x^*) \geq \frac{\mu_\omega}{2} \omega^2(x - x^*). \quad (21)$$

By convexity of  $f$ , we have  $\langle x - x^*, \nabla f(x) \rangle \geq f(x) - f(x^*)$ , and by definition of the Fenchel conjugate, we have

$$\omega(x - x^*) \cdot \omega_*(\nabla f(x)) \geq \langle x - x^*, \nabla f(x) \rangle \geq f(x) - f(x^*).$$

Hence by plugging (21), we obtain (20). ■

We now prove Theorem D.4 that is similar to Theorem 4.4. It states that in the case of the FW algorithm, the modified directional smoothness constant is bounded if the function is smooth, strongly convex and the set is strongly convex for any distance function  $\omega$ . It also provides an explicit upper bound on the modified directional smoothness constant. This bound implies that the convergence rate in Theorem D.2 is better conditioned than existing results [Garber et al. 2015].

**Theorem D.4 (Bounds on modified directional smoothness)** *Consider  $x_0 \in \mathcal{C}$  and a function  $f$ , smooth w.r.t. the distance function  $\omega$ , with constant  $L_\omega$ , strongly convex w.r.t. the distance function  $\omega$ , with constant  $\mu_\omega$ , and the set  $\mathcal{C}$ , strongly convex with constant  $\alpha_\omega$ . Let  $\delta(x) = x - v(x)$ ,  $v(x)$  being the FW corner. Then, the function  $f(x)$  is modified directionally smooth w.r.t. to  $\delta$ , with constant*

$$\tilde{\mathcal{L}}_{f,\delta}(x_0) \leq \frac{\kappa_\omega \sqrt{2} L_\omega}{\alpha_\omega \sqrt{\mu_\omega}} \frac{1}{\sqrt{f(x_0) - f^*}}. \quad (22)$$

**Proof.** Let  $h \in [0, 1]$ . With the smoothness of  $f$ , we have

$$f(x + h\delta(x)) \leq f(x) - h \langle -\nabla f(x), \delta(x) \rangle + \frac{h^2 L_\omega}{2} \omega(\delta(x))^2.$$

Recall that when  $\delta(x)$  is the Frank-Wolfe direction, we have that the Frank-Wolfe gap  $g(x)$  is equal to  $\langle -\nabla f(x), \delta(x) \rangle$ . Also, the scaling inequality for strongly convex sets (Lemma 3.6) implies that  $\omega(\delta(x))^2 \leq g(x)/(\alpha_\omega \omega^*(-\nabla f(x)))$ , so that

$$f(x + h\delta(x)) \leq f(x) - h \langle -\nabla f(x), \delta(x) \rangle + \frac{h^2 L_\omega}{2\alpha_\omega} \frac{g(x)}{\omega^*(-\nabla f(x))}.$$

<sup>1</sup>Note that Garber et al. [2015] use a different argument for the base case, bounding instead  $h_1$  with  $L \cdot \text{diam}(\mathcal{C})^2/2$ , using the Lipschitz smoothness of  $f$  (and this would become  $C_f/2$  in its affine invariant formulation with  $C_f$  as defined by Jaggi [2013]). However, we believe that  $h_0$  is usually smaller than  $C_f$  in applications, and in any case  $h_0$  appears from  $1/M^2$  for us, so using our different base case argument is more meaningful.

Now, it is easy to see from the definition of the dual distance  $\omega_*$  that it has the same bounded asymmetry constant as for  $\omega$ , and thus  $\omega^*(-\nabla f(x)) \geq \frac{1}{\kappa_\omega} \omega^*(\nabla f(x))$ . Thus we apply (20) to obtain:

$$f(x + h\delta(x)) \leq f(x) - hg(x) + \frac{h^2}{2} \frac{\kappa_w \sqrt{2} L_\omega}{\alpha_\omega \sqrt{\mu_\omega} \sqrt{f(x_0) - f^*}} \frac{\sqrt{f(x_0) - f^*}}{\sqrt{f(x) - f^*}} g(x),$$

which implies equation (22). ■

Theorem D.4 shows that the conditioning of convergence with the directional smoothness, which does not depend on any norm choice, in Theorem D.2 is better than conditioning of other analysis [Garber et al. 2015]. We now prove that the optimal constant of modified directional smoothness  $\tilde{L}_{f,\delta}$  is affine invariant, a result similar to Proposition 4.3 for the directional smoothness constant.

**Proposition D.5 (Affine Invariance of Modified Directional Smoothness)** *Consider  $\mathcal{C}$  a compact convex set and  $f$  a convex function on  $\mathcal{C}$  that is modified directionally smooth w.r.t.  $\delta(x)$  with constant  $\tilde{\mathcal{L}}_{f,\delta}(x_0)$  (with  $x_0 \in \mathcal{C}$ ). If for any  $x \in \mathcal{C}$ ,  $\delta(x)$  is affine covariant (e.g. the Frank-Wolfe direction  $\delta(x) \triangleq v(x) - x$ ), then the constant  $\tilde{\mathcal{L}}_{f,\delta}$  in (17) is affine invariant. In other words, for an invertible matrix  $B$ , let*

$$\tilde{f}(\cdot) \triangleq f(B\cdot), \quad \tilde{\delta}_{\tilde{\mathcal{C}}}(\cdot) \triangleq \delta_{B^{-1}\mathcal{C}}(\cdot),$$

then  $\tilde{\mathcal{L}}_{\tilde{f},\tilde{\delta}_{\tilde{\mathcal{C}}}}(x_0) = \tilde{\mathcal{L}}_{f,\delta}(y_0)$ , where  $y_0 \triangleq B^{-1}x_0$ .

**Proof.** Let  $y \in B^{-1} \cdot \mathcal{C}$ . Applying the definition of directional smoothness for  $f$  at  $By$ , we obtain

$$f(By + h\delta(By)) \leq f(By) + h\langle \nabla f(By), \delta(By) \rangle - \frac{\tilde{\mathcal{L}}_{f,\delta}(x_0)h^2}{2} \langle \nabla f(By), \delta(By) \rangle \sqrt{\frac{f(x_0) - f^*}{f(By) - f^*}}. \quad (23)$$

Similarly to Proposition 4.3, we have that  $\nabla \tilde{f}(y) = B^T \nabla f(By)$  and  $\delta(By) = B\tilde{\delta}_{\tilde{\mathcal{C}}}(y)$  so that

$$\langle \nabla f(By), \delta(By) \rangle = \langle \nabla f(By), B\tilde{\delta}_{\tilde{\mathcal{C}}}(y) \rangle = \langle B^T \nabla f(By), \tilde{\delta}_{\tilde{\mathcal{C}}}(y) \rangle = \langle \nabla \tilde{f}(y), \tilde{\delta}_{\tilde{\mathcal{C}}}(y) \rangle.$$

Hence (23) and  $\tilde{f}^* = f^*$ , implies that for any  $y \in B^{-1} \cdot \mathcal{C}$

$$\tilde{f}(y + h\tilde{\delta}_{\tilde{\mathcal{C}}}(y)) \leq \tilde{f}(y) + h\langle \nabla \tilde{f}(y), \tilde{\delta}_{\tilde{\mathcal{C}}}(y) \rangle - \frac{\tilde{\mathcal{L}}_{f,\delta}(x_0)h^2}{2} \langle \nabla \tilde{f}(y), \tilde{\delta}_{\tilde{\mathcal{C}}}(y) \rangle \sqrt{\frac{\tilde{f}(y_0) - \tilde{f}^*}{\tilde{f}(y) - \tilde{f}^*}}.$$

Hence,  $\tilde{f}$  is modified directionally smooth on  $\tilde{\mathcal{C}} \triangleq B^{-1} \cdot \mathcal{C}$  with respect to  $\tilde{\delta}_{\tilde{\mathcal{C}}}$  and  $\tilde{L}_{\tilde{f},\tilde{\delta}_{\tilde{\mathcal{C}}}}(y_0) \leq \tilde{\mathcal{L}}_{f,\delta}(x_0)$ . A similar reasoning concludes that the two constants are equal. ■

## E Related Work Details

Lacoste-Julien et al. 2013 propose an affine invariant analysis of the vanilla Frank-Wolfe algorithm when the unconstrained optimum  $x^*$  is in the relative interior of the constraint set  $\mathcal{C}$  and  $f$  is strongly convex. Hence, the analysis applies when the constraint set is a strongly convex set, and the quantity might be defined in our context. However, the affine invariant constant  $\mu_f^{(FW)}$  standing for the strong convexity of  $f$  is zero whenever the optimum is not in the relative interior of the constraint set  $\mathcal{C}$ . Indeed, Equation (3) from [Lacoste-Julien et al. 2013] define the following affine invariant quantity

$$\mu_f^{(FW)} \triangleq \inf_{\substack{x \in \mathcal{C} \setminus \{x^*\}, \gamma \in ]0,1] \\ \bar{s} = \bar{s}(x, x^*, \mathcal{C}) \\ y = x + \gamma(\bar{s} - x)}} \frac{2}{\gamma^2} [f(y) - f(x) - \langle \nabla f(x), y - x \rangle],$$

where  $\bar{s}(x, x^*, \mathcal{C}) = \text{ray}(x, x^*) \cap \partial\mathcal{C}$ . When  $x^* \notin \mathcal{C}$ , we have  $\mu_f^{(FW)} \leq 0$  since there are some point  $x \in \partial\mathcal{C}$  such that  $x \in \bar{s}(x, x^*, \mathcal{C})$ , and thus we can take  $\bar{s} = x$  in the inf, yielding  $y = x$  with  $\gamma > 0$ . This means that the above quantity cannot be easily generalized to the setting we studied in Theorem 4.4 where the unconstrained optimum is assumed to be *outside* of  $\mathcal{C}$ .