
Affine Invariant Analysis of Frank-Wolfe on Strongly Convex Sets

Thomas Kerdreux^{*1} Lewis Liu^{*2,3} Simon Lacoste Julien^{2,3,4,5} Damien Scieur^{*4,3}

Abstract

It is known that the Frank-Wolfe (FW) algorithm, which is affine covariant, enjoys faster convergence rates than $\mathcal{O}(1/K)$ 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. We show that our rates are better than any other known convergence rates of FW in this setting. 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 present similar performances than its affine invariant counterpart, despite using affine dependent norms in the step size’s computation.

1. Introduction

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

$$\min_{x \in \mathcal{C}} f(x), \quad \mathcal{C} \text{ convex and compact.} \quad (1)$$

FW algorithms 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 et al., 2018; Lacoste-Julien

^{*}Equal contribution ¹Zuse Institute, Berlin ²Département d’informatique et de recherche opérationnelle (DIRO), Université de Montréal ³Mila, Montréal ⁴Samsung SAIT AI Lab, Montréal ⁵Canada CIFAR AI Chair. Correspondence to: Thomas Kerdreux <thomaskerdreux@gmail.com>, Damien Scieur <damienscieur@gmail.com>.

Algorithm 1 Frank-Wolfe Algorithm

Input: $x_0 \in \mathcal{C}$.
1: **for** $k = 0, 1, \dots, K$ **do**
2: $v_k \in \operatorname{argmax}_{v \in \mathcal{C}} \langle -\nabla f(x_k), v - x_k \rangle$ \triangleright LMO
3: $\gamma_k = \operatorname{argmin}_{\gamma \in [0,1]} f(x_k + \gamma(v_k - x_k))$ \triangleright Line-search
4: $x_{k+1} = (1 - \gamma_k)x_k + \gamma_k v_k$ \triangleright Convex update
5: **end for**

et al., 2015; Courty et al., 2016; Paty & Cuturi, 2019; Luise et al., 2019; Combettes & Pokutta, 2021); however, many open questions remain in designing such optimal algorithmic schemes (*e.g.* (Braun et al., 2017; Kerdreux et al., 2018; Braun et al., 2019; Combettes & Pokutta, 2020; Carderera & Pokutta, 2020; Mortagy et al., 2020; Combettes et al., 2020; Bomze et al., 2021)) 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),

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

Definition 1.1 (Affine covariance) *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 enjoys a slow sublinear rate $\mathcal{O}(1/K)$ over general compact convex sets 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. Importantly, this analysis is better than any other known *best norm-dependent* analysis. (By *best norm-dependent analysis*, we refer to

the norm that minimizes the convergence rate of the theoretical analysis that depend on norms, see, e.g., (Lan, 2013, 3.13.)).

Definition 1.2 (norm-independence) *A quantity is norm-independent if it does not depend on the choice of a norm.*

Counterexample. *The condition number in optimization – the ratio between the smoothness and the strong convexity constants (Nesterov, 2013) – is norm-dependent. Therefore, algorithms whose rate depends on this condition number may be faster if the choice of the norm makes the condition number closer to 1.*

Example. *The curvature constant C_f (Jaggi, 2013) is defined by the ratio*

$$C_f \stackrel{\text{def}}{=} \sup_{\substack{x, v \in C \\ \gamma \in [0, 1] \\ y = x + \gamma(v - x)}} \frac{1}{\gamma^2} \left[f(y) - f(x) - \langle y - x; \nabla f(x) \rangle \right],$$

where C is a compact, convex set. Since this ratio does not involve any norm, it is therefore norm-independent.

Affine invariance and norm-independence are closely related, although they are quite different in nature. We discuss extensively their common points and differences in Appendix A. However, since the FW algorithm is affine invariant and norm-independent, its analysis should ideally also satisfy such properties.

Many works have then sought to find structural assumptions and algorithmic modifications that accelerate this sub-linear rate of $\mathcal{O}(1/K)$. The strong convexity of the set (or more generally uniform convexity, see (Kerdreux et al., 2021b;a)) 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 & Polyak, 1966; Demyanov & Rubinov, 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 & Hazan, 2015). However, to the best of our knowledge, there exists no norm-independent affine invariant analysis for these accelerated regimes.

In these “non affine invariant” analyses, structural assumptions like the L -smoothness (Definition 1.3) of f and the α -strong convexity of \mathcal{C} (Definition 1.4) lead to accelerated convergence rate of the Frank-Wolfe algorithm, but are typically conditioned on parameters L, α 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 1.3 *The function f is smooth over \mathcal{C} 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 1.4 *A set \mathcal{C} is α -strongly convex if, for any $(x, y) \in \mathcal{C}$, $\gamma \in [0, 1]$ and $\|z\| \leq 1$, we have*

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

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.

While the smoothness of a function is quite a standard assumption, the strong convexity of a set is a rather strong assumption. Nevertheless, strong convexity of sets are common in machine learning applications. We can cite, for instance, ℓ_p norms (common regularization in machine learning problems or action set in online learning) (Bubeck et al., 2018; Kerdreux et al., 2021c; Wang et al., 2021), matrix Schatten norms (Braverman et al., 2020), and matrix group norms (Kakade et al., 2012).

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 best accelerated convergence regimes of Frank-Wolfe on strongly convex sets (note, however, the recent proof of an accelerated asymptotic $\mathcal{O}(1/T^2)$ rate of vanilla Frank-Wolfe for specific scheduled step sizes (Bach, 2020)). 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 norm-independent, affine invariant rates for Frank-Wolfe 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.

Affine Invariant Analysis of Frank-Wolfe on Strongly Convex Sets

Related Work	\mathcal{C}	Str. cvx. f	x^*	Algo	Step size	Rate
Clarkson (2010)	Simplex	✗	Any	FW	Scheduled	$\mathcal{O}(1/K)$
Jaggi (2013)	Convex	✗	Any	FW	Scheduled	$\mathcal{O}(1/K)$
Lacoste-Julien & Jaggi (2013)	Any	✓	Interior	FW	Exact ls	Linear
Jaggi & Lacoste-Julien (2015)	Polytope	✓	Any	Corr. FW	Exact ls	Linear
Gutman & Pena (2020)						
Our work	Strongly cvx	✗	$\nabla f(x^*) \neq 0$	FW	Backtracking ls	Linear
	Strongly cvx	✓	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 & Jaggi, 2013; Jaggi & Lacoste-Julien, 2015), the authors consider “generalized geometric strong convexity” (see their Eq. 39), an affine invariant measure of (generalized) strong convexity, while (Gutman & Pena, 2020) consider strongly convex functions relative to a pair (\mathcal{C}, ω) where ω 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).

Contributions. In this paper, **1)** we conduct affine invariant analysis of the Frank-Wolfe Algorithm 1, when the function f is smooth and the set \mathcal{C} is strongly convex. Our affine invariant conditioning is better than any norm-dependent analysis. Additionally, we point out that there is likely a positive gap between our constant and the optimal norm-dependent bound, given that ours are not restricted to a choice of same norms for different parameters in the bound. In specific, we 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 w.r.t. a specified direction δ* . Based on this definition, **2)** we propose an affine invariant backtracking line-search for finding the optimal step size, which achieves the best of two worlds in theory and practice. Finally, **3)** we show that existing backtracking line-search methods, which use a specific norm, converges surprisingly to the optimal norm-independent, affine invariant step size. This implies that affine dependent and affine invariant backtracking techniques perform similarly.

Outline. In Section 2, we motivate the need for norm-independent 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 best affine invariant analy-

sis. For instance, corrective variants of Frank-Wolfe exhibit (affine invariant) linear convergence rates when the constraint set is a polytope (Lacoste-Julien & Jaggi, 2013; Jaggi & Lacoste-Julien, 2015) and the objective function is (generally) strongly convex. See Table 1 for a review of all affine invariant analyses of Frank-Wolfe algorithms.

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 (Kerdreux, 2020), *e.g.* Hölderian Error Bounds on f (Kerdreux et al., 2019; Xu & Yang, 2018; Rinaldi & Zeffiro, 2020) or local uniform convexity of \mathcal{C} (Kerdreux et al., 2021b).

Our affine invariant analyses introduce constants seeking to characterize structural properties without a specific choice of norm, even the best (inaccessible) one (see Appendix A for an in-depth discussion). This has been the basis for works extending the accelerated convergence analysis to non-smooth or non-strongly convex functions (Pena, 2019; Gutman & Pena, 2020), which then explore new structural assumptions on f .

Gap between affine invariant and best-norm analysis. We point out that, in general, affine invariance does not imply optimality. For instance, even if designing norms that produce affine invariant rates is possible (d’Aspremont et al., 2018), this does not imply that such rates will match the result of our norm-independent, affine invariant analy-

sis. In this paper, we show that our affine invariant constant is always better than norm-dependent ones, even after taking the best norm.

Furthermore, it is still an open question if there is a gap between affine invariant rates and the best-norm rate comprising of norm-dependent parameters (such as smoothness, strong convexity and the lower bound of gradient norms). We believe this may be the case, since, the best-norm rate implicitly impose the same norm on all the parameters in the bound, while our affine invariant constant is free of such constraints.

To conclude, we highlight that characterizing the gap between affine invariant and best norm analysis is an interesting and challenging problem in the literature. See Appendix A for more details and examples on this problem.

Notations. For a norm $\|\cdot\|$, we write $\|\cdot\|_*$ its dual norm. Our ambient space is \mathbb{R}^d .

2. Norm-Dependent Analysis of FW

It is known that when the function is *smooth* (Definition 1.3), the set is *strongly-convex* (Definition 1.4) and the gradient is lower bounded $\|\nabla f(x)\|_* \geq c > 0$ over the constraint set (i.e., the constraints are active), the Frank-Wolfe algorithm 1 converges linearly (Levitin & Polyak, 1966; Demyanov & Rubinov, 1970; Dunn, 1979), at rate (with $h_k \stackrel{\text{def}}{=} f(x_k) - f^*$)

$$h_k \leq \left(\max \left\{ \frac{1}{2}, 1 - \frac{c\alpha}{4L} \right\} \right)^k h_0. \quad (5)$$

see (Garber & Hazan, 2015; Kerdreux et al., 2021b) for more details, and we recall these results in Appendix B, in respectively Lemma B.1 and Corollary B.2. 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 , α , 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 , α depend on “affine variant” norms.

Example 2.1 Consider the projection of $\bar{x} : \|\bar{x}\|_2 > 1$,

$$\min_x f(x) \stackrel{\text{def}}{=} \frac{1}{2} \|x - \bar{x}\|_2^2 \quad \text{such that } \|x\|_2^2 \leq 1.$$

In such case, we have that $L = 1$, $\alpha = 1$ and $c = \|\bar{x}\|_2 - 1$ (L , α and c are defined according to the ℓ_2 norm, see proof

in Appendix B.2). However, if we transform the problem into $\min_y f(By)$, the new constants become

$$L = \sigma_{\max}(B), \quad \alpha = \frac{\sigma_{\min}(B)}{\sigma_{\max}(B)}, \quad c = \sigma_{\max}(B)(\|\bar{x}\|_2 - 1).$$

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

$$\begin{aligned} f(x_k) - f^* &\leq \left(1 - \frac{\|\bar{x}\|_2 - 1}{4} \right)^k (f(x_0) - f^*), \\ f(By_k) - f^* &\leq \left(1 - \frac{\|\bar{x}\|_2 - 1}{4} \kappa^{-1}(B) \right)^k (f(x_0) - f^*), \end{aligned}$$

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 (dues to their affine covariance), therefore the upper bound will not be representative of the true rate of convergence of FW.

Remark. The constants, and therefore the rate, can be improved if we change the norm $\|\cdot\|_2$ into $\|\cdot\|_{2,B^{-1}}$. However, it is usually very hard or impossible to guess what norm will be the best for a specific problem. This is not a problem for FW with exact line-search, as no norm is required. However, in the case of (Garber & Hazan, 2015), the step size (or backtracking line-search) strategy uses L , and therefore the rate depends directly on the choice of the norm. Moreover, even if we choose the gauge of the Euclidean ball to measure the function smoothness and the set strong-convexity (becoming, in this case, invariant to affine reparametrization of our problem, see Appendix A), we do not know how to guarantee it was the optimal choice for this specific problem.

When the optimum is in the relative interior of any compact set \mathcal{C} , FW converges linearly when f is strongly convex (Guélat & Marcotte, 1986; Lacoste-Julien & Jaggi, 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 & Rubinov, 1970). Our paper hence focuses on extending the analysis where the unconstrained optimum is outside the constraint set (Demyanov & Rubinov, 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 & Hazan, 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 E, we provide an affine invariant analysis of (Garber & Hazan, 2015).

3. Scaling Inequalities on Strongly Convex Sets

All proofs of Frank-Wolfe methods on strongly convex sets leverage the same property. The *scaling inequality* crucially relates the Frank-Wolfe gap with $\|x_t - v_t\|^2$, see e.g. (Kerdreux et al., 2021b, Lemma 2.1.). The scaling inequality is an equivalent characterization of strong convexity of \mathcal{C} (Goncharov & Ivanov, 2017, Theorem 2.1.), but we recall here only the implication that we will need, see (Kerdreux et al., 2021a) for a review of useful properties of uniformly convex sets in machine learning. Importantly, the scaling inequalities motivate the new structural assumptions we present in Section 4 and Appendix E.

Lemma 3.1 (Norm Scaling Inequality) *Assume \mathcal{C} is α -strongly convex w.r.t. a norm $\|\cdot\|$. 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*

$$\langle \phi, v_\phi - x \rangle \geq \frac{\alpha}{2} \|\phi\|_* \|v_\phi - x\|^2. \quad (6)$$

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

$$\langle -\nabla f(x_k); v_k - x_k \rangle \geq \frac{\alpha}{2} \|\nabla f(x_k)\|_* \|v_k - x_k\|^2.$$

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,

$$\gamma x + (1-\gamma)v_\phi + \frac{\alpha}{2}\gamma(1-\gamma)\|x - v_\phi\|^2 z \in \mathcal{C} \quad \forall z : \|z\| \leq 1.$$

Then, by optimality of v_ϕ on \mathcal{C} ,

$$\langle \phi; v_\phi \rangle \geq \langle \phi; \gamma x + (1-\gamma)v_\phi + \frac{\alpha}{2}\gamma(1-\gamma)\|x - v_\phi\|^2 z \rangle$$

After simplification,

$$\langle \phi; v_\phi - x \rangle \geq \frac{\alpha}{2}(1-\gamma)\|x - v_\phi\|^2 \langle \phi; z \rangle.$$

With $\gamma \rightarrow 0$, and after maximizing over z , we obtain by definition of $\|\cdot\|_*$,

$$\langle \phi; v_\phi - x \rangle \geq \frac{\alpha}{2} \|x - v_\phi\|^2 \|\phi\|_*,$$

which holds in particular when $\phi = -\nabla f(x)$. ■

These scaling inequalities can take other forms as in the following corollary.

Corollary 3.2 *Assume \mathcal{C} is α -strongly convex w.r.t. $\|\cdot\|$. Consider $(d_1, d_2) \in \mathbb{R}^d$ s.t. $\min\{\|d_1\|_*, \|d_2\|_*\} > c > 0$ and let $(x_1, x_2) \in \partial\mathcal{C}$, s.t. $d_i \in N_{\mathcal{C}}(x_i)$ for $i = 1, 2$. Then*

$$\|x_1 - x_2\| \leq \|d_1 - d_2\|_*/(\alpha c).$$

Proof. By applying successively Lemma 3.1, we obtain

$$\begin{aligned} \langle d_1; x_1 - x_2 \rangle &\geq \alpha/2 \|d_1\|_* \|x_1 - x_2\|^2 \\ \langle d_2; x_2 - x_1 \rangle &\geq \alpha/2 \|d_2\|_* \|x_1 - x_2\|^2. \end{aligned}$$

We then obtain $\langle d_1 - d_2; x_1 - x_2 \rangle \geq \alpha c \|x_1 - x_2\|^2$. Finally, by definition of the dual norm, we conclude that $\alpha c \|x_1 - x_2\| \leq \|d_1 - d_2\|_*$. ■

This Corollary provides new insights on (Lan, 2013, Algorithm 4). Indeed, it implies that when the set \mathcal{C} is strongly convex and $\inf_{x \in \mathcal{C}} \|\nabla f(x)\|_* > c > 0$, then the strong condition on the Linear Minimization Oracle (Lan, 2013, Equation (4.4.)) is satisfied with $\rho = 1$ and hence PA-CndG (Lan, 2013, Algorithm 4) converges in $\mathcal{O}(1/K^2)$ (Lan, 2013, Corollary 1).

PA-CndG is a Frank-Wolfe type algorithm with the oblivious step-sizes $\frac{2}{k+2}$, hence affine co-variant. Note, however, that the $\mathcal{O}(1/K^2)$ accelerated convergence rate is achieved under the same structural assumption that ensures linear convergence of Frank-Wolfe in (Levitin & Polyak, 1966), which on the other hand, require exact line-search or problem-dependent step-sizes.

4. Directional Smoothness

Analyses of Frank-Wolfe algorithm on strongly convex sets 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 δ . 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 $x + h\delta(x) \in \mathcal{C}$,*

$$\begin{aligned} f(x + h\delta(x)) &\leq f(x) - h \langle -\nabla f(x), \delta(x) \rangle \\ &\quad + \frac{\mathcal{L}_{f,\delta} h^2}{2} \langle -\nabla f(x), \delta(x) \rangle. \end{aligned} \quad (7)$$

The rationale of Definition 4.1 is to replace the norm in the usual smoothness condition (Definition 1.3) 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 (7) 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 ,

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) - h\|\nabla f(x)\|_2^2 + \frac{Lh^2}{2}\|\nabla f(x)\|_2^2 \\ &= f(x) + h\left(\frac{Lh}{2} - 1\right)\|\nabla f(x)\|_2^2. \end{aligned}$$

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 $\delta(x) \stackrel{\text{def}}{=} v(x) - x$), then $\mathcal{L}_{f,\delta}$ in (7) is invariant to an affine transformation of the constraint set (proof in Appendix B.3).

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 norm $\|\cdot\|$.

Theorem 4.4 (Directional Smoothness of FW) Consider the function f , smooth w.r.t. the norm $\|\cdot\|$, with constant $L_{\|\cdot\|}$, and the set \mathcal{C} , strongly convex with constant $\alpha_{\|\cdot\|}$. Let $\delta(x) = v(x) - x$, $v(x)$ being the FW vertex

$$v(x) \stackrel{\text{def}}{=} \operatorname{argmin}_{v \in \mathcal{C}} \langle \nabla f(x), v \rangle. \quad (8)$$

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

$$\mathcal{L}_{f,\delta} \leq 2 \frac{L_{\|\cdot\|}}{c_{\|\cdot\|} \alpha_{\|\cdot\|}}. \quad (9)$$

Proof. See Appendix B.4 for the proof. ■

5. Affine Invariant Linear Rates

With the directional smoothness constant $\mathcal{L}_{f,\delta}$ (affine invariant when δ 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 & Polyak, 1966; Demyanov & Rubinov, 1970; Dunn, 1979).

Theorem 5.1 (Affine Invariant Linear Rates) Assume f is a convex function and directionally smooth with direction function δ 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 vertex (8), converges linearly, at rate

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

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

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

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

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

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 corollary shows that the directional constant in Theorem 5.1 is bounded by (9) w.r.t. the norm $\|\cdot\|$ that gives the best ratio.

Corollary 5.2 Write Ω the set of norms in \mathbb{R}^d . Then, the rate of convergence using directional smoothness is at least better than the previously known, norm-dependent rate,

$$1 - \frac{1}{2\mathcal{L}_{f,\delta}} \leq 1 - \frac{1}{4 \min_{\|\cdot\| \in \Omega} \frac{L_{\|\cdot\|}}{c_{\|\cdot\|} \alpha_{\|\cdot\|}}},$$

where $L_{\|\cdot\|}$ is the smoothness constant of the function f , $\alpha_{\|\cdot\|}$ the strong convexity of the set \mathcal{C} and $\inf_{x \in \mathcal{C}} \|\nabla f(x)\|_* = c_{\|\cdot\|} > 0$.

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

In Appendix E, we provide an affine invariant analysis without restriction on the position of the optimum, i.e. the $\mathcal{O}(1/K^2)$ analysis in (Garber & Hazan, 2015). We define (Definition E.1) a similar property to the directional smoothness that additionally accounts for the strong convexity of f . We choose to present the affine invariant analysis for the linear convergence in the main body of the paper 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 FW. The significant advantage of these constants is that, like FW, they are independent of any norm. However, the optimal step size of FW 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 technique finds automatically an estimate of the directional smoothness that satisfies

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

at the cost of one additional function evaluation per iteration. It is known that such backtracking technique is, in the worst case, two times slower than FW with the optimal, affine invariant stepsize.

Algorithm 2 Affine invariant backtracking

Input: FW vertex 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 (7),

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

- 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 D.

Algorithm 3 in Appendix D 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 ℓ_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 is at worst four times smaller than the theoretical affine invariant stepsize.

Proposition 7.1 *Let f be directionally smooth, and let $L(x) = \frac{\mathcal{L}_{f,\delta} \langle \nabla f(x), \delta(x) \rangle}{\|\delta(x)\|_2^2}$. Assume $L(x)$ locally approximately constant, i.e., there exists k_{\min}, k_{\max} such that, for $L_{loc} = \max_i L(x_i)$,*

$$\frac{L_{loc}}{2} < L(x_k) \leq L_{loc}, \quad k \in [k_{\min}, k_{\max}].$$

In this case, the norm-dependent backtracking line-search Algorithm 3 finds

$$L_k < 2L_{loc}, \quad k = \left\lceil k_{\min} + \log_2 \frac{L_{k_{\min}}}{L_{loc}} \right\rceil, \dots, k_{\max},$$

and its step size $(\gamma_)_k$ satisfies*

$$\min \left\{ 1, \frac{1}{4\mathcal{L}_{f,\delta}} \right\} \leq (\gamma_*)_k.$$

Proof. See Appendix B.5 for the full proof.

Proof sketch. The constant L_{loc} can be seen as the local Lipschitz constant. Indeed, if we write the upper bound given by the directional smoothness, we have

$$\begin{aligned} f(x) + h \langle \nabla f(x), \delta(x) \rangle + \frac{h^2}{2} \mathcal{L}_{f,\delta} \langle \nabla f(x_k), \delta(x_k) \rangle \\ = f(x) + h \langle \nabla f(x), \delta(x) \rangle + L(x) \frac{h^2}{2} \|\delta(x)\|_2^2, \end{aligned}$$

where the right-hand-side corresponds to the definition of smoothness (3) at $y = x + \delta(x)$ with a variable constant $L(x)$. The parameter $L(x)$ can thus be seen as a “local Lipschitz constant”. If $L(x)$ remains approximately constant, the backtracking line-search will eventually find an estimation $L_k \leq 2L_{loc}$. Therefore, with the norm-dependent backtracking line-search, the step size will be at worst 4 times smaller than the one of the affine invariant fixed-step strategy. ■

Therefore, the optimal step size from the backtracking line-search with the ℓ_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,

$$\min_{x \in \mathcal{C}} \frac{1}{n} \sum_{i=1}^n l(a_i^T x, y_i), \quad (10)$$

where l is the quadratic or the logistic loss. Here we adopt the ℓ_2 -ball, defined as

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

Specifically, we compare our affine invariant backtracking method in Algorithm 2 against the naive FW Algorithm 1 with step size $1/L$ (Demjanov & Rubinov, 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 x^* satisfies $\|x^*\|_2 = 1.1R$, and the initial iterate $x_0 = \mathbf{0}$. 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 ℓ_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 for smooth functions 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, 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. Research reported in this paper was also partially supported through the Research Campus Modal funded by the German Federal Ministry of Education and Research (fund numbers 05M14ZAM,05M20ZBM) as well as the Deutsche Forschungsgemeinschaft (DFG) through the DFG Cluster of Excellence MATH+.

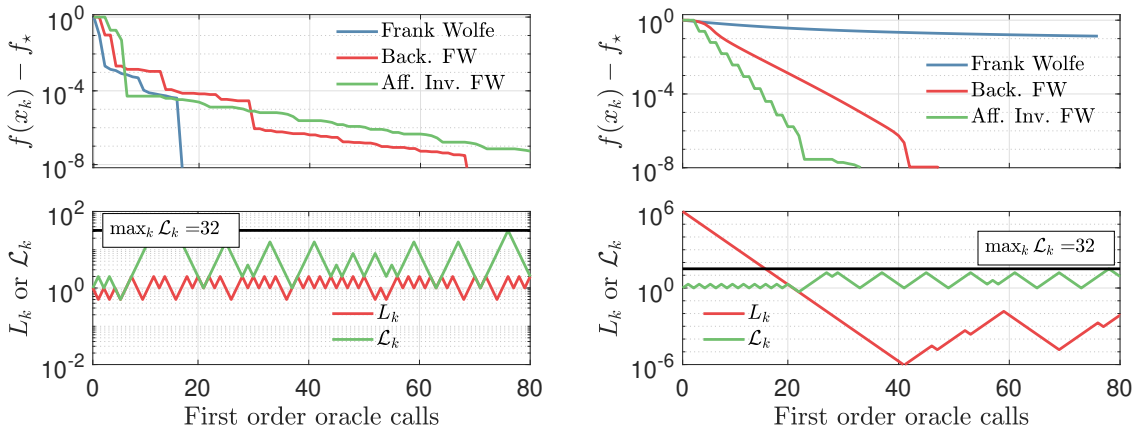


Figure 1. Comparison of FW variants on the projection problem. Left: $B = I$, Right: $\kappa(B) = 10^6$. The top row is the gap $f_k - f^*$, and the bottom row corresponds to the estimation of the directional-smoothness constant \mathcal{L}_k or the smoothness constant L_k , where the black line report 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.

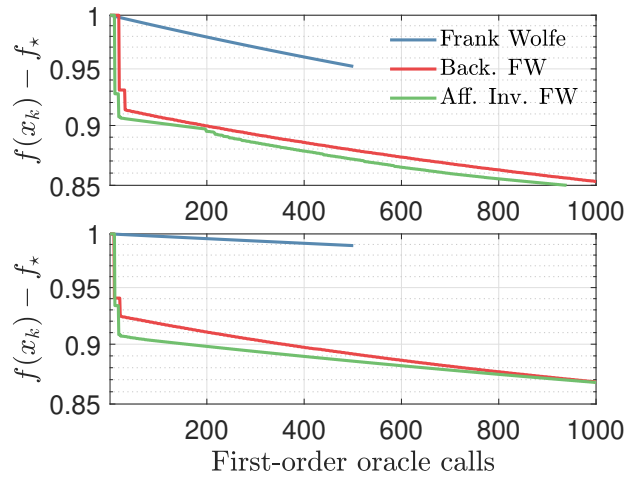


Figure 2. Classification problem on Madelon dataset, with (Top) Quadratic loss and (Bottom) Logistic loss.

References

- Alayrac, J.-B., Bojanowski, P., Agrawal, N., Sivic, J., Laptev, I., and Lacoste-Julien, S. Unsupervised learning from narrated instruction videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4575–4583, 2016.
- Bach, F. On the effectiveness of Richardson extrapolation in machine learning. *arXiv preprint arXiv:2002.02835*, 2020.
- Bauschke, H. H., Bolte, J., and Teboulle, M. A descent lemma beyond lipschitz gradient continuity: first-order methods revisited and applications. *Mathematics of Operations Research*, 42(2):330–348, 2017.
- Bojanowski, P., Lajugie, R., Bach, F., Laptev, I., Ponce, J., Schmid, C., and Sivic, J. Weakly supervised action labeling in videos under ordering constraints. In *European Conference on Computer Vision*. Springer, 2014.
- Bomze, I. M., Rinaldi, F., and Zeffiro, D. Fast cluster detection in networks by first-order optimization. *arXiv preprint arXiv:2103.15907*, 2021.
- Braun, G., Pokutta, S., and Zink, D. Lazifying conditional gradient algorithms. *Proceedings of ICML*, 2017.
- Braun, G., Pokutta, S., Tu, D., and Wright, S. Blended conditional gradients. In *International Conference on Machine Learning*, pp. 735–743. PMLR, 2019.
- Braverman, V., Krauthgamer, R., Krishnan, A., and Sinoff, R. Schatten norms in matrix streams: Hello sparsity, goodbye dimension. In *International Conference on Machine Learning*, pp. 1100–1110. PMLR, 2020.
- Bubeck, S., Cohen, M., and Li, Y. Sparsity, variance and curvature in multi-armed bandits. In *Algorithmic Learning Theory*, pp. 111–127. PMLR, 2018.
- Carderera, A. and Pokutta, S. Second-order conditional gradients. *arXiv preprint arXiv:2002.08907*, 2020.
- Clarkson, K. Coresets, sparse greedy approximation, and the Frank-Wolfe algorithm. *ACM Transactions on Algorithms (TALG)*, 6(4):63, 2010.
- Combettes, C. and Pokutta, S. Boosting Frank-Wolfe by chasing gradients. In *International Conference on Machine Learning*, pp. 2111–2121. PMLR, 2020.
- Combettes, C. W. and Pokutta, S. Complexity of linear minimization and projection on some sets. *arXiv:2101.10040*, 2021.
- Combettes, C. W., Spiegel, C., and Pokutta, S. Projection-free adaptive gradients for large-scale optimization. *arXiv:2009.14114*, 2020.
- Courty, N., Flamary, R., Tuia, D., and Rakotomamonjy, A. Optimal transport for domain adaptation. *IEEE transactions on pattern analysis and machine intelligence*, 39(9):1853–1865, 2016.
- d’Aspremont, A., Guzman, C., and Jaggi, M. Optimal affine-invariant smooth minimization algorithms. *SIAM Journal on Optimization*, 28(3):2384–2405, 2018.
- Demyanov, V. F. and Rubinov, A. M. Approximate methods in optimization problems. *Modern Analytic and Computational Methods in Science and Mathematics*, 1970.
- Dunn, J. C. Rates of convergence for conditional gradient algorithms near singular and nonsingular extremals. *SIAM Journal on Control and Optimization*, 17(2):187–211, 1979.
- Frank, M., Wolfe, P., et al. An algorithm for quadratic programming. *Naval research logistics quarterly*, 3(1-2): 95–110, 1956.
- Garber, D. and Hazan, E. Faster rates for the Frank-Wolfe method over strongly-convex sets. In *32nd International Conference on Machine Learning, ICML 2015*, 2015.
- Goncharov, V. V. and Ivanov, G. E. Strong and weak convexity of closed sets in a Hilbert space. In *Operations research, engineering, and cyber security*, pp. 259–297. Springer, 2017.
- Guélat, J. and Marcotte, P. Some comments on Wolfe’s ‘away step’. *Mathematical Programming*, 1986.
- Gutman, D. H. and Pena, J. F. The condition number of a function relative to a set. *Mathematical Programming*, pp. 1–40, 2020.
- Guyon, I., Li, J., Mader, T., Pletscher, P. A., Schneider, G., and Uhr, M. Competitive baseline methods set new standards for the nips 2003 feature selection benchmark. *Pattern recognition letters*, 28(12):1438–1444, 2007.
- Jaggi, M. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In *Proceedings of the 30th international conference on machine learning*, number CONF, pp. 427–435, 2013.
- Jaggi, M. and Lacoste-Julien, S. On the global linear convergence of Frank-Wolfe optimization variants. *Advances in Neural Information Processing Systems*, 28, 2015.
- Journée, M., Nesterov, Y., Richtárik, P., and Sepulchre, R. Generalized power method for sparse principal component analysis. *Journal of Machine Learning Research*, 11(2), 2010.

- Kakade, S. M., Shalev-Shwartz, S., and Tewari, A. Regularization techniques for learning with matrices. *The Journal of Machine Learning Research*, 13(1):1865–1890, 2012.
- Kerdreux, T. *Accelerating conditional gradient methods*. PhD thesis, Université Paris sciences et lettres, 2020.
- Kerdreux, T., Pedregosa, F., and d’Aspremont, A. Frank-Wolfe with subsampling oracle. In *International Conference on Machine Learning*, pp. 2591–2600. PMLR, 2018.
- Kerdreux, T., d’Aspremont, A., and Pokutta, S. Restarting Frank-Wolfe. In *The 22nd International Conference on Artificial Intelligence and Statistics*. PMLR, 2019.
- Kerdreux, T., d’Aspremont, A., and Pokutta, S. Local and global uniform convexity conditions. *arXiv preprint arXiv:2102.05134*, 2021a.
- Kerdreux, T., d’Aspremont, A., and Pokutta, S. Projection-free optimization on uniformly convex sets. In *International Conference on Artificial Intelligence and Statistics*, pp. 19–27. PMLR, 2021b.
- Kerdreux, T., Roux, C., d’Aspremont, A., and Pokutta, S. Linear bandits on uniformly convex sets. *arXiv preprint arXiv:2103.05907*, 2021c.
- Lacoste-Julien, S. and Jaggi, M. An affine invariant linear convergence analysis for Frank-Wolfe algorithms. *arXiv preprint arXiv:1312.7864*, 2013.
- Lacoste-Julien, S., Lindsten, F., and Bach, F. Sequential kernel herding: Frank-wolfe optimization for particle filtering. In *Artificial Intelligence and Statistics*, pp. 544–552. PMLR, 2015.
- Lan, G. The complexity of large-scale convex programming under a linear optimization oracle. *arXiv preprint arXiv:1309.5550*, 2013.
- Levitin, E. S. and Polyak, B. T. Constrained minimization methods. *USSR Computational mathematics and mathematical physics*, 6(5):1–50, 1966.
- Lu, H., Freund, R. M., and Nesterov, Y. Relatively smooth convex optimization by first-order methods, and applications. *SIAM Journal on Optimization*, 28(1):333–354, 2018.
- Luise, G., Salzo, S., Pontil, M., and Ciliberto, C. Sinkhorn barycenters with free support via Frank-Wolfe algorithm. In *Advances in Neural Information Processing Systems*, pp. 9318–9329, 2019.
- Miech, A., Laptev, I., and Sivic, J. Learning a text-video embedding from incomplete and heterogeneous data. *arXiv preprint arXiv:1804.02516*, 2018.
- Molinaro, M. Curvature of feasible sets in offline and online optimization. *arXiv:2002.03213*, 2020.
- Mortagy, H., Gupta, S., and Pokutta, S. Walking in the shadow: A new perspective on descent directions for constrained minimization. *Advances in Neural Information Processing Systems*, 33, 2020.
- Nesterov, Y. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013.
- Paty, F.-P. and Cuturi, M. Subspace robust wasserstein distances. In *International Conference on Machine Learning*, pp. 5072–5081. PMLR, 2019.
- Pedregosa, F., Negiar, G., Askari, A., and Jaggi, M. Linearly convergent Frank-Wolfe with backtracking line-search. In *International Conference on Artificial Intelligence and Statistics*, pp. 1–10. PMLR, 2020.
- Pena, J. Generalized conditional subgradient and generalized mirror descent: duality, convergence, and symmetry. *arXiv preprint arXiv:1903.00459*, 2019.
- Peyre, J., Sivic, J., Laptev, I., and Schmid, C. Weakly-supervised learning of visual relations. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 5179–5188, 2017.
- Rector-Brooks, J., Wang, J.-K., and Mozafari, B. Revisiting projection-free optimization for strongly convex constraint sets. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 1576–1583, 2019.
- Rinaldi, F. and Zeffiro, D. A unifying framework for the analysis of projection-free first-order methods under a sufficient slope condition. *arXiv:2008.09781*, 2020.
- Rockafellar, R. T. *Convex analysis*. Princeton university press, 1970.
- Seguin, G., Bojanowski, P., Lajugie, R., and Laptev, I. Instance-level video segmentation from object tracks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- Wang, H., Yang, X., and Deng, X. A hybrid first-order method for nonconvex lp-ball constrained optimization. *arXiv preprint arXiv:2104.04400*, 2021.
- Xu, Y. and Yang, T. Frank-Wolfe method is automatically adaptive to error bound condition. *arXiv:1810.04765*, 2018.