



# Frank–Wolfe and friends: a journey into projection-free first-order optimization methods

Immanuel M. Bomze<sup>1</sup> · Francesco Rinaldi<sup>2</sup> · Damiano Zeffiro<sup>2</sup>

Received: 1 July 2021 / Revised: 29 July 2021 / Accepted: 18 August 2021 /  
Published online: 6 September 2021  
© The Author(s) 2021

## Abstract

Invented some 65 years ago in a seminal paper by Marguerite Straus-Frank and Philip Wolfe, the Frank–Wolfe method recently enjoys a remarkable revival, fuelled by the need of fast and reliable first-order optimization methods in Data Science and other relevant application areas. This review tries to explain the success of this approach by illustrating versatility and applicability in a wide range of contexts, combined with an account on recent progress in variants, improving on both the speed and efficiency of this surprisingly simple principle of first-order optimization.

**Keywords** First-order methods · Projection-free methods · Structured optimization · Conditional gradient · Sparse optimization

**Mathematics Subject Classification** 90C06 · 90C25 · 90C30

## 1 Introduction

In their seminal work (Frank and Wolfe 1956), Marguerite Straus-Frank and Philip Wolfe introduced a first-order algorithm for the minimization of convex quadratic objectives over polytopes, now known as Frank–Wolfe (FW) method. The main idea of the method is simple: to generate a sequence of feasible iterates by moving at every step towards a minimizer of a linearized objective, the so-called FW vertex.

---

✉ Immanuel M. Bomze  
immanuel.bomze@univie.ac.at

Francesco Rinaldi  
rinaldi@math.unipd.it

Damiano Zeffiro  
damiano.zeffiro@math.unipd.it

<sup>1</sup> ISOR, VCOR & ds:UniVie, Universität Wien, Vienna, Austria

<sup>2</sup> Dipartimento di Matematica “Tullio Levi-Civita”, Università di Padova, Padua, Italy

Subsequent works, partly motivated by applications in optimal control theory (see Dunn (1979) for references), generalized the method to smooth (possibly non-convex) optimization over closed subsets of Banach spaces admitting a linear minimization oracle (see Demyanov and Rubinov 1970; Dunn and Harshbarger 1978).

Furthermore, while the  $\mathcal{O}(1/k)$  rate in the original article was proved to be optimal when the solution lies on the boundary of the feasible set (Canon and Cullum 1968), improved rates were given in a variety of different settings. In Levitin and Polyak (1966) and Demyanov and Rubinov (1970), a linear convergence rate was proved over strongly convex domains assuming a lower bound on the gradient norm, a result then extended in Dunn (1979) under more general gradient inequalities. In Guélat and Marcotte (1986), linear convergence of the method was proved for strongly convex objectives with the minimum obtained in the relative interior of the feasible set.

The slow convergence behaviour for objectives with solution on the boundary motivated the introduction of several variants, the most popular being Wolfe's away step (Wolfe 1970). Wolfe's idea was to move away from bad vertices, in case a step of the FW method moving towards good vertices did not lead to sufficient improvement on the objective. This idea was successfully applied in several network equilibrium problems, where linear minimization can be achieved by solving a min-cost flow problem (see Fukushima 1984 and references therein). In Guélat and Marcotte (1986), some ideas already sketched by Wolfe were formalized to prove linear convergence of the Wolfe's away step method and identification of the face containing the solution in finite time, under some suitable strict complementarity assumptions.

In recent years, the FW method has regained popularity thanks to its ability to handle the structured constraints appearing in machine learning and data science applications efficiently. Examples include LASSO, SVM training, matrix completion, minimum enclosing ball, density mixture estimation, cluster detection, to name just a few (see Sect. 3 for further details).

One of the main features of the FW algorithm is its ability to naturally identify sparse and structured (approximate) solutions. For instance, if the optimization domain is the simplex, then after  $k$  steps the cardinality of the support of the last iterate generated by the method is at most  $k + 1$ . Most importantly, in this setting every vertex added to the support at every iteration must be the best possible in some sense, a property that connects the method with many greedy optimization schemes (Clarkson 2010). This makes the FW method pretty efficient on the abovementioned problem class. Indeed, the combination of structured solutions with often noisy data makes the sparse approximations found by the method possibly more desirable than high precision solutions generated by a faster converging approach. In some cases, like in cluster detection (see, e.g., Bomze 1997), finding the support of the solution is actually enough to solve the problem independently from the precision achieved.

Another important feature is that the linear minimization used in the method is often cheaper than the projections required by projected-gradient methods. It is important to notice that, even when these two operations have the same complexity, constants defining the related bounds can differ significantly (see Combettes and Pokutta 2021 for some examples and tests). When dealing with large scale problems, the FW method hence has a much smaller per-iteration cost with respect to projected-gradient methods. For this reason, FW methods fall into the category of *projection-free methods*

(Lan 2020). Furthermore, the method can be used to approximately solve quadratic subproblems in accelerated schemes, an approach usually referred to as conditional gradient sliding (see, e.g., Carderera and Pokutta 2020; Lan and Zhou 2016).

## 1.1 Organisation of the paper

The present review is not intended to provide an exhaustive literature survey, but rather as an advanced tutorial demonstrating versatility and power of this approach. The article is structured as follows: in Sect. 2, we introduce the classic FW method, together with a general scheme for all the methods we consider. In Sect. 3, we present applications from classic optimization to more recent machine learning problems. In Sect. 4, we review some important stepsizes for first order methods. In Sect. 5, we discuss the main theoretical results about the FW method and the most popular variants, including the  $\mathcal{O}(1/k)$  convergence rate for convex objectives, affine invariance, the sparse approximation property, and support identification. In Sect. 6 we illustrate some recent improvements on the  $\mathcal{O}(1/k)$  convergence rate. Finally, in Sect. 7 we present recent FW variants fitting different optimization frameworks, in particular block coordinate, distributed, accelerated, and trace norm optimization. We highlight that all the proofs reported in the paper are either seminal, or simplified versions of proofs reported in published papers, and we believe they might give some useful technical insights to the interested reader.

## 1.2 Notation

For any integers  $a$  and  $b$ , denote by  $[a : b] = \{x \text{ integer} : a \leq x \leq b\}$  the integer range between them. For a set  $V$ , the power set  $2^V$  denotes the system all subsets of  $V$ , whereas for any positive integer  $s \in \mathbb{N}$  we set  $\binom{V}{s} := \{S \in 2^V : |S| = s\}$ , with  $|S|$  denoting the number of elements in  $S$ . Matrices are denoted by capital sans-serif letters (e.g., the zero matrix  $\mathbf{O}$ , or the  $n \times n$  identity matrix  $\mathbf{I}_n$  with columns  $\mathbf{e}_i$  the length of which should be clear from the context). The all-ones vector is  $\mathbf{e} := \sum_i \mathbf{e}_i \in \mathbb{R}^n$ . Generally, vectors are always denoted by boldface sans-serif letters  $\mathbf{x}$ , and their transpose by  $\mathbf{x}^\top$ . The Euclidean norm of  $\mathbf{x}$  is then  $\|\mathbf{x}\| := \sqrt{\mathbf{x}^\top \mathbf{x}}$  whereas the general  $p$ -norm is denoted by  $\|\mathbf{x}\|_p$  for any  $p \geq 1$  (so  $\|\mathbf{x}\|_2 = \|\mathbf{x}\|$ ). By contrast, the so-called zero-norm simply counts the number of nonzero entries:

$$\|\mathbf{x}\|_0 := |\{i \in [1:n] : x_i \neq 0\}|.$$

For a vector  $\mathbf{d}$  we denote as  $\widehat{\mathbf{d}} := \frac{1}{\|\mathbf{d}\|} \mathbf{d}$  its normalization, with the convention  $\widehat{\mathbf{d}} = \mathbf{o}$  if  $\mathbf{d} = \mathbf{o}$ . Here  $\mathbf{o}$  denotes the zero vector. In context of symmetric matrices, “psd” abbreviates “positive-semidefinite”.

## 2 Problem and general scheme

We consider the following problem:

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

where, unless specified otherwise,  $C$  is a convex and compact (i.e. bounded and closed) subset of  $\mathbb{R}^n$  and  $f$  is a differentiable function having Lipschitz continuous gradient with constant  $L > 0$ :

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\| \quad \text{for all } \{x, y\} \subset C.$$

This is a central property required in the analysis of first-order methods. Such a property indeed implies (and for a convex function is equivalent to) the so-called Descent Lemma (see, e.g., Bertsekas 2015, Proposition 6.1.2), which provides a quadratic upper approximation to the function  $f$ . Throughout the article, we denote by  $x^*$  a (global) solution to (1) and use the symbol  $f^* := f(x^*)$  as a shorthand for the corresponding optimal value.

The general scheme of the first-order methods we consider for problem (1), reported in Algorithm 1, is based upon a set  $F(x, g)$  of directions feasible at  $x$  using first-order local information on  $f$  around  $x$ , in the smooth case  $g = \nabla f(x)$ . From this set, a particular  $d \in F(x, g)$  is selected, with the maximal stepsize  $\alpha^{\max}$  possibly dependent from auxiliary information available to the method (at iteration  $k$ , we thus write  $\alpha_k^{\max}$ ), and not always equal to the maximal feasible stepsize.

---

### Algorithm 1 First-order method

---

- 1 Choose a point  $x_0 \in C$
  - 2 For  $k = 0, \dots$
  - 3     If  $x_k$  satisfies some specific condition, then STOP
  - 4     Choose  $d_k \in F(x_k, \nabla f(x_k))$
  - 5     Set  $x_{k+1} = x_k + \alpha_k d_k$ , with  $\alpha_k \in (0, \alpha_k^{\max}]$  a suitably chosen stepsize
  - 6 End for
- 

### 2.1 The classical Frank–Wolfe method

The classical FW method for minimization of a smooth objective  $f$  generates a sequence of feasible points  $\{x_k\}$  following the scheme of Algorithm 2. At the iteration  $k$  it moves toward a vertex i.e., an extreme point, of the feasible set minimizing the scalar product with the current gradient  $\nabla f(x_k)$ . It therefore makes use of a linear minimization oracle (LMO) for the feasible set  $C$

$$\text{LMO}_C(g) \in \arg \min_{z \in C} g^T z, \quad (2)$$

defining the descent direction as

$$\mathbf{d}_k = \mathbf{d}_k^{FW} := \mathbf{s}_k - \mathbf{x}_k, \quad \mathbf{s}_k \in \text{LMO}_C(\nabla f(x_k)). \quad (3)$$

In particular, the update at step 6 in Algorithm 2 can be written as

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k(\mathbf{s}_k - \mathbf{x}_k) = \alpha_k \mathbf{s}_k + (1 - \alpha_k)\mathbf{x}_k \quad (4)$$

Since  $\alpha_k \in [0, 1]$ , by induction  $\mathbf{x}_{k+1}$  can be written as a convex combination of elements in the set  $S_{k+1} := \{\mathbf{x}_0\} \cup \{\mathbf{s}_i\}_{0 \leq i \leq k}$ . When  $C = \text{conv}(A)$  for a set  $A$  of points with some common property, usually called “elementary atoms”, if  $\mathbf{x}_0 \in A$  then  $\mathbf{x}_k$  can be written as a convex combination of  $k + 1$  elements in  $A$ . Note that due to Caratheodory’s theorem, we can even limit the number of occurring atoms to  $\min\{k, n\} + 1$ . In the rest of the paper the primal gap at iteration  $k$  is defined as  $h_k = f(\mathbf{x}_k) - f^*$ .

---

### Algorithm 2 Frank–Wolfe method

---

- 1 Choose a point  $\mathbf{x}_0 \in C$
  - 2 For  $k = 0, \dots$
  - 3     If  $\mathbf{x}_k$  satisfies some specific condition, then STOP
  - 4     Compute  $\mathbf{s}_k \in \text{LMO}_C(\nabla f(x_k))$
  - 5     Set  $\mathbf{d}_k^{FW} = \mathbf{s}_k - \mathbf{x}_k$
  - 6     Set  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k^{FW}$ , with  $\alpha_k \in (0, 1]$  a suitably chosen stepsize
  - 7 End for
- 

## 3 Examples

FW methods and variants are a natural choice for constrained optimization on convex sets admitting a linear minimization oracle significantly faster than computing a projection. We present here in particular the traffic assignment problem, submodular optimization, LASSO problem, matrix completion, adversarial attacks, minimum enclosing ball, SVM training, maximal clique search in graphs, sparse optimization.

### 3.1 Traffic assignment

Finding a traffic pattern satisfying the equilibrium conditions in a transportation network is a classic problem in optimization that dates back to Wardrop’s paper (Wardrop 1952). Let  $\mathcal{G}$  be a network with set of nodes  $[1 : n]$ . Let  $\{D(i, j)\}_{i \neq j}$  be demand coefficients, modeling the amount of goods with destination  $j$  and origin  $i$ . For any  $i, j$  with  $i \neq j$  let furthermore  $f_{ij} : \mathbb{R} \rightarrow \mathbb{R}$  be the non-linear (convex) cost functions, and  $x_{ij}^s$  be the flow on link  $(i, j)$  with destination  $s$ . The traffic assignment problem can be modeled as the following non-linear *multicommodity network* problem (Fukushima

1984):

$$\min \left\{ \sum_{i,j} f_{ij} \left( \sum_s x_{ij}^s \right) : \sum_i x_{\ell i}^s - \sum_j x_{j\ell}^s = D(\ell, s), \text{ all } \ell \neq s, x_{ij}^s \geq 0 \right\}. \quad (5)$$

Then the linearized optimization subproblem necessary to compute the FW vertex takes the form

$$\min \left\{ \sum_s \sum_{i,j} c_{ij} x_{ij}^s : \sum_i x_{\ell i}^s - \sum_j x_{j\ell}^s = D(\ell, s), \ell \neq s, x_{ij}^s \geq 0 \right\} \quad (6)$$

and can be split in  $n$  shortest paths subproblems, each of the form

$$\min \left\{ \sum_{i,j} c_{ij} x_{ij}^s : \sum_i x_{\ell i}^s - \sum_j x_{j\ell}^s = D(\ell, s), \ell \neq s, x_{ij}^s \geq 0 \right\} \quad (7)$$

for a fixed  $s \in [1:n]$ , with  $c_{ij}$  the first-order derivative of  $f_{ij}$  (see Fukushima 1984 for further details). A number of FW variants were proposed in the literature for efficiently handling this kind of problems (see, e.g., Bertsekas 2015; Fukushima 1984; LeBlanc et al. 1975; Mitradjieva and Lindberg 2013; Weintraub et al. 1985 and references therein for further details). Some of those variants represent a good (if not the best) choice when low or medium precision is required in the solution of the problem (Perederieieva et al. 2015).

In the more recent work Joulin et al. (2014) a FW variant also solving a shortest path subproblem at each iteration was applied to image and video co-localization.

### 3.2 Submodular optimization

Given a finite set  $V$ , a function  $r : 2^V \rightarrow \mathbb{R}$  is said to be submodular if for every  $A, B \subset V$

$$r(A) + r(B) \geq r(A \cup B) + r(A \cap B). \quad (8)$$

As is common practice in the optimization literature (see e.g. Bach 2013, Section 2.1), here we always assume  $s(\emptyset) = 0$ . A number of machine learning problems, including image segmentation and sensor placement, can be cast as minimization of a submodular function (see, e.g., Bach 2013; Chakrabarty et al. 2014 and references therein for further details):

$$\min_{A \subset V} r(A). \quad (9)$$

Submodular optimization can also be seen as a more general way to relate combinatorial problems to convexity, for example for structured sparsity (Bach 2013; Jaggi

2013). By a theorem from Fujishige (1980), problem (9) can be in turn reduced to an minimum norm point problem over the base polytope

$$B(r) = \left\{ \mathbf{s} \in \mathbb{R}^V : \sum_{a \in A} s_a \leq r(A) \text{ for all } A \subseteq V, \sum_{a \in V} s_a = r(V) \right\}. \quad (10)$$

For this polytope, linear optimization can be achieved with a simple greedy algorithm. More precisely, consider the LP

$$\max_{\mathbf{s} \in B(r)} \mathbf{w}^\top \mathbf{s}.$$

Then if the objective vector  $\mathbf{w}$  has a negative component, the problem is clearly unbounded. Otherwise, a solution to the LP can be obtained by ordering  $\mathbf{w}$  in decreasing manner as  $w_{j_1} \geq w_{j_2} \geq \dots \geq w_{j_n}$ , and setting

$$s_{j_k} := r(\{j_1, \dots, j_k\}) - r(\{j_1, \dots, j_{k-1}\}), \quad (11)$$

for  $k \in [1:n]$ . We thus have a LMO with a  $\mathcal{O}(n \log n)$  cost. This is the reason why FW variants are widely used in the context of submodular optimization; further details can be found in, e.g., Bach (2013), Jaggi (2013).

### 3.3 LASSO problem

The LASSO, proposed by Tibshirani in 1996 (Tibshirani 1996), is a popular tool for sparse linear regression. Given the training set

$$T = \{(r_i, b_i) \in \mathbb{R}^n \times \mathbb{R} : i \in [1:m]\},$$

where  $\mathbf{r}_i^\top$  are the rows of an  $m \times n$  matrix  $\mathbf{A}$ , the goal is finding a sparse linear model (i.e., a model with a small number of non-zero parameters) describing the data. This problem is strictly connected with the Basis Pursuit Denoising (BPD) problem in signal analysis (see, e.g., Chen et al. 2001). In this case, given a discrete-time input signal  $b$ , and a *dictionary*

$$\{\mathbf{a}_j \in \mathbb{R}^m : j \in [1:n]\}$$

of elementary discrete-time signals, usually called atoms (here  $\mathbf{a}_j$  are the columns of a matrix  $\mathbf{A}$ ), the goal is finding a sparse linear combination of the atoms that *approximate* the real signal. From a purely formal point of view, LASSO and BPD problems are equivalent, and both can be formulated as follows:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) &:= \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 \\ \text{s.t.} \quad &\|\mathbf{x}\|_1 \leq \tau, \end{aligned} \quad (12)$$

where the parameter  $\tau$  controls the amount of shrinkage that is applied to the model (related to sparsity, i.e., the number of nonzero components in  $x$ ). The feasible set is

$$C = \{x \in \mathbb{R}^n : \|x\|_1 \leq \tau\} = \text{conv}\{\pm\tau e_i : i \in [1:n]\}.$$

Thus we have the following LMO in this case:

$$\text{LMO}_C(\nabla f(x_k)) = \text{sign}(-\nabla_{i_k} f(x_k)) \cdot \tau e_{i_k},$$

with  $i_k \in \arg \max_i |\nabla_i f(x_k)|$ . It is easy to see that the FW per-iteration cost is then  $\mathcal{O}(n)$ . The peculiar structure of the problem makes FW variants well-suited for its solution. This is the reason why LASSO/BPD problems were considered in a number of FW-related papers (see, e.g., Jaggi 2011, 2013; Lacoste-Julien and Jaggi 2015; Locatello et al. 2017).

### 3.4 Matrix completion

Matrix completion is a widely studied problem that comes up in many areas of science and engineering, including collaborative filtering, machine learning, control, remote sensing, and computer vision (just to name a few; see also Candès and Recht (2009) and references therein). The goal is to retrieve a low rank matrix  $X \in \mathbb{R}^{n_1 \times n_2}$  from a sparse set of observed matrix entries  $\{U_{ij}\}_{(i,j) \in J}$  with  $J \subset [1:n_1] \times [1:n_2]$ . Thus the problem can be formulated as follows (Freund et al. 2017):

$$\begin{aligned} \min_{X \in \mathbb{R}^{n_1 \times n_2}} f(X) &:= \sum_{(i,j) \in J} (X_{ij} - U_{ij})^2 \\ \text{s.t.} \quad \text{rank}(X) &\leq \delta, \end{aligned} \tag{13}$$

where the function  $f$  is given by the squared loss over the observed entries of the matrix and  $\delta > 0$  is a parameter representing the assumed belief about the rank of the reconstructed matrix we want to get in the end. In practice, the low rank constraint is relaxed with a nuclear norm ball constraint, where we recall that the nuclear norm  $\|X\|_*$  of a matrix  $X$  is equal the sum of its singular values. Thus we get the following convex optimization problem:

$$\begin{aligned} \min_{X \in \mathbb{R}^{n_1 \times n_2}} \sum_{(i,j) \in J} (X_{ij} - U_{ij})^2 \\ \text{s.t.} \quad \|X\|_* \leq \delta. \end{aligned} \tag{14}$$

The feasible set is the convex hull of rank-one matrices:

$$\begin{aligned} C &= \{X \in \mathbb{R}^{n_1 \times n_2} : \|X\|_* \leq \delta\} \\ &= \text{conv}\{\delta uv^T : u \in \mathbb{R}^{n_1}, v \in \mathbb{R}^{n_2}, \|u\| = \|v\| = 1\}. \end{aligned}$$

If we indicate with  $A_J$  the matrix that coincides with  $A$  on the indices  $J$  and is zero otherwise, then we can write  $\nabla f(X) = 2(X - U)_J$ . Thus we have the following LMO

in this case:

$$\text{LMO}_C(\nabla f(X_k)) \in \arg \min\{\text{tr}(\nabla f(X_k)^\top X) : \|X\|_* \leq \delta\}, \quad (15)$$

which boils down to computing the gradient, and the rank-one matrix  $\delta u_1 v_1^\top$ , with  $u_1, v_1$  right and left singular vectors corresponding to the top singular value of  $-\nabla f(X_k)$ . Consequently, the FW method at a given iteration approximately reconstructs the target matrix as a sparse combination of rank-1 matrices. Furthermore, as the gradient matrix is sparse (it only has  $|J|$  non-zero entries) storage and approximate singular vector computations can be performed much more efficiently than for dense matrices<sup>1</sup>. A number of FW variants has hence been proposed in the literature for solving this problem (see, e.g., Freund et al. 2017; Jaggi 2011, 2013).

### 3.5 Adversarial attacks in machine learning

Adversarial examples are maliciously perturbed inputs designed to mislead a properly trained learning machine at test time. An *adversarial attack* hence consists in taking a correctly classified data point  $x_0$  and slightly modifying it to create a new data point that leads the considered model to misclassification (see, e.g., Carlini and Wagner 2017; Chen et al. 2017; Goodfellow et al. 2014 for further details). A possible formulation of the problem (see, e.g., Chen et al. 2020; Goodfellow et al. 2014) is given by the so called *maximum allowable  $\ell_p$ -norm attack* that is,

$$\begin{aligned} \min_{x \in \mathbb{R}^n} & f(x_0 + x) \\ \text{s.t.} & \|x\|_p \leq \varepsilon, \end{aligned} \quad (16)$$

where  $f$  is a suitably chosen attack loss function,  $x_0$  is a correctly classified data point,  $x$  represents the additive noise/perturbation,  $\varepsilon > 0$  denotes the magnitude of the attack, and  $p \geq 1$ . It is easy to see that the LMO has a cost  $\mathcal{O}(n)$ . If  $x_0$  is a feature vector of a dog image correctly classified by our learning machine, our adversarial attack hence suitably perturbs the feature vector (using the noise vector  $x$ ), thus getting a new feature vector  $x_0 + x$  classified, e.g., as a cat. In case a target adversarial class is specified by the attacker, we have a *targeted attack*. In some scenarios, the goal may not be to push  $x_0$  to a specific target class, but rather push it away from its original class. In this case we have a so called *untargeted attack*. The attack loss function  $f$  will hence be chosen depending on the kind of attack we aim to perform over the considered model. Due to its specific structure, problem (16) can be nicely handled by means of tailored FW variants. Some FW frameworks for adversarial attacks were recently described in, e.g., Chen et al. (2020), Kazemi et al. (2021), Sahu and Kar (2020).

<sup>1</sup> Details related to the LMO cost can be found in, e.g., Jaggi (2013).

### 3.6 Minimum enclosing ball

Given a set of points  $P = \{p_1, \dots, p_n\} \subset \mathbb{R}^d$ , the minimum enclosing ball problem (MEB, see, e.g., Clarkson 2010; Yıldırım 2008) consists in finding the smallest ball containing  $P$ . Such a problem models numerous important applications in clustering, nearest neighbor search, data classification, machine learning, facility location, collision detection, and computer graphics, to name just a few. We refer the reader to Kumar et al. (2003) and the references therein for further details. Denoting by  $c \in \mathbb{R}^d$  the center and by  $\sqrt{\gamma}$  (with  $\gamma \geq 0$ ) the radius of the ball, a convex quadratic formulation for this problem is

$$\min_{(c, \gamma) \in \mathbb{R}^d \times \mathbb{R}} \gamma \tag{17}$$

$$s.t. \|p_i - c\|^2 \leq \gamma, \quad \text{all } i \in [1:n]. \tag{18}$$

This problem can be formulated via Lagrangian duality as a convex *Standard Quadratic Optimization Problem* (StQP, see, e.g. Bomze and de Klerk 2002)

$$\min \{x^T A^T A x - b^T x : x \in \Delta_{n-1}\} \tag{19}$$

with  $A = [p_1, \dots, p_n]$  and  $b^T = [p_1^T p_1, \dots, p_n^T p_n]$ . The feasible set is the standard simplex

$$\Delta_{n-1} := \{x \in \mathbb{R}_+^n : e^T x = 1\} = \text{conv}\{e_i : i \in [1:n]\},$$

and the LMO is defined as follows:

$$\text{LMO}_{\Delta_{n-1}}(\nabla f(x_k)) = e_{i_k},$$

with  $i_k \in \arg \min_i \nabla_i f(x_k)$ . It is easy to see that cost per iteration is  $\mathcal{O}(n)$ . When applied to (19), the FW method can find an  $\varepsilon$ -cluster in  $\mathcal{O}(\frac{1}{\varepsilon})$ , where an  $\varepsilon$ -cluster is a subset  $P'$  of  $P$  such that the MEB of  $P'$  dilated by  $1 + \varepsilon$  contains  $P$  (Clarkson 2010). The set  $P'$  is given by the atoms in  $P$  selected by the LMO in the first  $\mathcal{O}(\frac{1}{\varepsilon})$  iterations. Further details related to the connections between FW methods and MEB problems can be found in, e.g., Ahıpařaoglu et al. (2008), Ahıpařaoglu and Todd (2013), Clarkson (2010) and references therein.

### 3.7 Training linear Support Vector Machines

*Support Vector Machines (SVMs)* represent a very important class of machine learning tools (see, e.g., Vapnik 2013 for further details). Given a labeled set of data points, usually called *training set*:

$$TS = \{(p_i, y_i), p_i \in \mathbb{R}^d, y_i \in \{-1, 1\}, i = 1, \dots, n\},$$

the linear SVM training problem consists in finding a linear classifier  $w \in \mathbb{R}^d$  such that the label  $y_i$  can be deduced with the “highest possible confidence” from  $w^\top p_i$ . A convex quadratic formulation for this problem is the following Clarkson (2010):

$$\begin{aligned} \min_{w \in \mathbb{R}^d, \rho \in \mathbb{R}} \quad & \rho + \frac{\|w\|^2}{2} \\ \text{s.t.} \quad & \rho + y_i w^\top p_i \geq 0, \quad \text{all } i \in [1:n], \end{aligned} \quad (20)$$

where the slack variable  $\rho$  stands for the negative margin and we can have  $\rho < 0$  if and only if there exists an exact linear classifier, i.e.  $w$  such that  $w^\top p_i = \text{sign}(y_i)$ . The dual of (20) is again an StQP:

$$\min \{x^\top A^\top A x : x \in \Delta_{n-1}\} \quad (21)$$

with  $A = [y_1 p_1, \dots, y_n p_n]$ . Notice that problem (21) is equivalent to an MNP problem on  $\text{conv}\{y_i p_i : i \in [1:n]\}$ , see Sect. 7.2 below. Some FW variants (like, e.g., the Pairwise Frank–Wolfe) are closely related to classical working set algorithms, such as the SMO algorithm used to train SVMs (Lacoste-Julien and Jaggi 2015). Further details on FW methods for SVM training problems can be found in, e.g., Clarkson (2010), Jaggi (2011).

### 3.8 Finding maximal cliques in graphs

In the context of network analysis the clique model, dating back at least to the work of Luce and Perry (1949) about social networks, refers to subsets with every two elements in a direct relationship. The problem of finding maximal cliques has numerous applications in domains including telecommunication networks, biochemistry, financial networks, and scheduling (see, e.g., Bomze et al. 1999; Wu and Hao 2015). Let  $G = (V, E)$  be a simple undirected graph with  $V$  and  $E$  set of vertices and edges, respectively. A clique in  $G$  is a subset  $C \subseteq V$  such that  $(i, j) \in E$  for each  $(i, j) \in C$ , with  $i \neq j$ . The goal in finding a clique  $C$  such that it is maximal (i.e., it is not contained in any strictly larger clique). This corresponds to find a local solution to the following equivalent (this time non-convex) StQP (see, e.g., Bomze 1997; Bomze et al. 1999; Hungerford and Rinaldi 2019 for further details):

$$\max \left\{ x^\top A_G x + \frac{1}{2} \|x\|^2 : x \in \Delta_{n-1} \right\} \quad (22)$$

where  $A_G$  is the adjacency matrix of  $G$ . Due to the peculiar structure of the problem, FW methods can be fruitfully used to find maximal cliques (see, e.g., Hungerford and Rinaldi 2019).

### 3.9 Finding sparse points in a set

Given a non-empty polyhedron  $P \subset \mathbb{R}^n$ , the goal is finding a sparse point  $x \in P$  (i.e., a point with as many zero components as possible). This sparse optimization

problem can be used to model a number of real-world applications in fields like, e.g., machine learning, pattern recognition and signal processing (see Rinaldi et al. 2010 and references therein). Ideally, what we would like to get is an optimal solution for the following problem:

$$\min \{ \|x\|_0 : x \in P \}. \quad (23)$$

Since the zero norm is non-smooth, a standard procedure is to replace the original formulation (23) with an equivalent concave optimization problem of the form:

$$\min \left\{ \sum_{i=1}^n \phi(y_i) : x \in P, -y \leq x \leq y \right\}, \quad (24)$$

where  $\phi : [0, +\infty[ \rightarrow \mathbb{R}$  is a suitably chosen smooth concave univariate function bounded from below, like, e.g.,

$$\phi(t) = (1 - e^{-\alpha t}),$$

with  $\alpha$  a large enough positive parameter (see, e.g., Mangasarian 1996; Rinaldi et al. 2010 for further details). The LMO in this case gives a vertex solution for the linear programming problem:

$$\min \{ c_k^T y : x \in P, -y \leq x \leq y \},$$

with  $(c_k)_i$  the first-order derivative of  $\phi$  calculated in  $(y_k)_i$ . Variants of the unit-stepsize FW method have been proposed in the literature (see, e.g., Mangasarian 1996; Rinaldi et al. 2010) to tackle the smooth equivalent formulation (24).

## 4 Stepsizes

Popular rules for determining the stepsize are:

- unit stepsize:

$$\alpha_k = 1,$$

mainly used when the problem has a concave objective function. Finite convergence can be proved, under suitable assumptions, both for the unit-stepsize FW and some of its variants described in the literature (see, e.g., Rinaldi et al. 2010 for further details).

- *diminishing stepsize*:

$$\alpha_k = \frac{2}{k+2}, \quad (25)$$

mainly used for the classic FW (see, e.g., Freund and Grigas 2016; Jaggi 2013).

– *exact line search*:

$$\alpha_k = \min_{\alpha \in [0, \alpha_k^{\max}]} \arg \min \varphi(\alpha) \quad \text{with } \varphi(\alpha) := f(\mathbf{x}_k + \alpha \mathbf{d}_k), \tag{26}$$

where we pick the smallest minimizer of the function  $\varphi$  for the sake of being well-defined even in rare cases of ties (see, e.g., Bomze et al. 2020; Lacoste-Julien and Jaggi 2015).

– *Armijo line search*: the method iteratively shrinks the step size in order to guarantee a sufficient reduction of the objective function. It represents a good way to replace exact line search in cases when it gets too costly. In practice, we fix parameters  $\delta \in (0, 1)$  and  $\gamma \in (0, \frac{1}{2})$ , then try steps  $\alpha = \delta^m \alpha_k^{\max}$  with  $m \in \{0, 1, 2, \dots\}$  until the sufficient decrease inequality

$$f(\mathbf{x}_k + \alpha \mathbf{d}_k) \leq f(\mathbf{x}_k) + \gamma \alpha \nabla f(\mathbf{x}_k)^\top \mathbf{d}_k \tag{27}$$

holds, and set  $\alpha_k = \alpha$  (see, e.g., Bomze et al. 2019 and references therein).

– *Lipschitz constant dependent step size*:

$$\alpha_k = \alpha_k(L) := \min \left\{ -\frac{\nabla f(\mathbf{x}_k)^\top \mathbf{d}_k}{L \|\mathbf{d}_k\|^2}, \alpha_k^{\max} \right\}, \tag{28}$$

with  $L$  the Lipschitz constant of  $\nabla f$  (see, e.g., Bomze et al. 2020; Pedregosa et al. 2020).

The Lipschitz constant dependent step size can be seen as the minimizer of the quadratic model  $m_k(\cdot; L)$  overestimating  $f$  along the line  $\mathbf{x}_k + \alpha \mathbf{d}_k$ :

$$m_k(\alpha; L) = f(\mathbf{x}_k) + \alpha \nabla f(\mathbf{x}_k)^\top \mathbf{d}_k + \frac{L\alpha^2}{2} \|\mathbf{d}_k\|^2 \geq f(\mathbf{x}_k + \alpha \mathbf{d}_k), \tag{29}$$

where the inequality follows by the standard Descent Lemma.

In case  $L$  is unknown, it is even possible to approximate  $L$  using a backtracking line search (see, e.g., Kerdreux et al. 2020; Pedregosa et al. 2020).

We now report a lower bound for the improvement on the objective obtained with the stepsize (28), often used in the convergence analysis.

**Lemma 1** *If  $\alpha_k$  is given by (28) and  $\alpha_k < \alpha_k^{\max}$  then*

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \frac{1}{2L} (\nabla f(\mathbf{x}_k)^\top \widehat{\mathbf{d}}_k)^2. \tag{30}$$

**Proof** We have

$$\begin{aligned} f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) &\leq f(\mathbf{x}_k) + \alpha_k \nabla f(\mathbf{x}_k)^\top \mathbf{d}_k + \frac{L\alpha_k^2}{2} \|\mathbf{d}_k\|^2 \\ &= f(\mathbf{x}_k) - \frac{(\nabla f(\mathbf{x}_k)^\top \mathbf{d}_k)^2}{2L \|\mathbf{d}_k\|^2} = f(\mathbf{x}_k) - \frac{1}{2L} (\nabla f(\mathbf{x}_k)^\top \widehat{\mathbf{d}}_k)^2, \end{aligned} \tag{31}$$

where we used the standard Descent Lemma in the inequality. □

## 5 Properties of the FW method and its variants

### 5.1 The FW gap

A key parameter often used as a measure of convergence is the FW gap

$$G(\mathbf{x}) = \max_{\mathbf{s} \in C} -\nabla f(\mathbf{x})^\top (\mathbf{s} - \mathbf{x}), \quad (32)$$

which is always nonnegative and equal to 0 only in first order stationary points. This gap is, by definition, readily available during the algorithm. If  $f$  is convex, using that  $\nabla f(\mathbf{x})$  is a subgradient we obtain

$$G(\mathbf{x}) \geq -\nabla f(\mathbf{x})^\top (\mathbf{x}^* - \mathbf{x}) \geq f(\mathbf{x}) - f^*, \quad (33)$$

so that  $G(\mathbf{x})$  is an upper bound on the optimality gap at  $\mathbf{x}$ . Furthermore,  $G(\mathbf{x})$  is a special case of the Fenchel duality gap (Lacoste-Julien et al. 2013).

If  $C = \Delta_{n-1}$  is the simplex, then  $G$  is related to the Wolfe dual as defined in Clarkson (2010). Indeed, this variant of Wolfe's dual reads

$$\begin{aligned} \max \quad & f(\mathbf{x}) + \lambda(\mathbf{e}^\top \mathbf{x} - 1) - \mathbf{u}^\top \mathbf{x} \\ \text{s.t.} \quad & \nabla_i f(\mathbf{x}) - u_i + \lambda = 0, \quad i \in [1:n], \\ & (\mathbf{x}, \mathbf{u}, \lambda) \in \mathbb{R}^n \times \mathbb{R}_+^n \times \mathbb{R} \end{aligned} \quad (34)$$

and for a fixed  $\mathbf{x} \in \mathbb{R}^n$ , the optimal values of  $(\mathbf{u}, \lambda)$  are

$$\lambda_{\mathbf{x}} = -\min_j \nabla_j f(\mathbf{x}), \quad u_i(\mathbf{x}) := \nabla_i f(\mathbf{x}) - \min_j \nabla_j f(\mathbf{x}) \geq 0.$$

Performing maximization in problem (34) iteratively, first for  $(\mathbf{u}, \lambda)$  and then for  $\mathbf{x}$ , this implies that (34) is equivalent to

$$\begin{aligned} \max_{\mathbf{x} \in \mathbb{R}^n} [f(\mathbf{x}) + \lambda_{\mathbf{x}}(\mathbf{e}^\top \mathbf{x} - 1) - \mathbf{u}(\mathbf{x})^\top \mathbf{x}] \\ = \max_{\mathbf{x} \in \mathbb{R}^n} [f(\mathbf{x}) - \max_j (\mathbf{e}_j - \mathbf{x})^\top \nabla f(\mathbf{x})] = \max_{\mathbf{x} \in \mathbb{R}^n} [f(\mathbf{x}) - G(\mathbf{x})]. \end{aligned} \quad (35)$$

Furthermore, since Slater's condition is satisfied, strong duality holds by Slater's theorem (Boyd et al. 2004), resulting in  $G(\mathbf{x}^*) = 0$  for every solution  $\mathbf{x}^*$  of the primal problem.

The FW gap is related to several other measures of convergence (see e.g. Lan 2020, Section 7.5.1). First, consider the projected gradient

$$\tilde{\mathbf{g}}_k := \pi_C(\mathbf{x}_k - \nabla f(\mathbf{x}_k)) - \mathbf{x}_k. \quad (36)$$

with  $\pi_B$  the projection on a convex and closed subset  $B \subseteq \mathbb{R}^n$ . We have  $\|\tilde{\mathbf{g}}_k\| = 0$  if and only if  $\mathbf{x}_k$  is stationary, with

$$\begin{aligned} \|\tilde{\mathbf{g}}_k\|^2 &= \tilde{\mathbf{g}}_k^\top \tilde{\mathbf{g}}_k \leq \tilde{\mathbf{g}}_k^\top [(\mathbf{x}_k - \nabla f(\mathbf{x}_k)) - \pi_C(\mathbf{x}_k - \nabla f(\mathbf{x}_k))] + \tilde{\mathbf{g}}_k^\top \tilde{\mathbf{g}}_k \\ &= -\tilde{\mathbf{g}}_k^\top \nabla f(\mathbf{x}_k) = -(\pi_C(\mathbf{x}_k - \nabla f(\mathbf{x}_k)) - \mathbf{x}_k)^\top \nabla f(\mathbf{x}_k) \\ &\leq \max_{\mathbf{y} \in C} -(\mathbf{y} - \mathbf{x}_k)^\top \nabla f(\mathbf{x}_k) = G(\mathbf{x}_k), \end{aligned} \tag{37}$$

where we used  $[\mathbf{y} - \pi_C(\mathbf{x})]^\top [\mathbf{x} - \pi_C(\mathbf{x})] \leq 0$  in the first inequality, with  $\mathbf{x} = \mathbf{x}_k - \nabla f(\mathbf{x}_k)$  and  $\mathbf{y} = \mathbf{x}_k$ .

Let now  $N_C(x)$  denote the normal cone to  $C$  at a point  $\mathbf{x} \in C$ :

$$N_C(\mathbf{x}) := \{\mathbf{r} \in \mathbb{R}^n : \mathbf{r}^\top (\mathbf{y} - \mathbf{x}) \leq 0 \text{ for all } \mathbf{y} \in C\}. \tag{38}$$

First-order stationarity conditions are equivalent to  $-\nabla f(\mathbf{x}) \in N_C(\mathbf{x})$ , or

$$\text{dist}(N_C(\mathbf{x}), -\nabla f(\mathbf{x})) = \|-\nabla f(\mathbf{x}) - \pi_{N_C(\mathbf{x})}(-\nabla f(\mathbf{x}))\| = 0.$$

The FW gap provides a lower bound on the distance from the normal cone  $\text{dist}(N_C(\mathbf{x}), -\nabla f(\mathbf{x}))$ , inflated by the diameter  $D > 0$  of  $C$ , as follows:

$$\begin{aligned} G(\mathbf{x}_k) &= -(\mathbf{s}_k - \mathbf{x}_k)^\top \nabla f(\mathbf{x}_k) \\ &= (\mathbf{s}_k - \mathbf{x}_k)^\top [\pi_{N_C(\mathbf{x}_k)}(-\nabla f(\mathbf{x}_k)) - (\pi_{N_C(\mathbf{x}_k)}(-\nabla f(\mathbf{x}_k)) + \nabla f(\mathbf{x}_k))] \\ &\leq \|\mathbf{s}_k - \mathbf{x}_k\| \|\pi_{N_C(\mathbf{x}_k)}(-\nabla f(\mathbf{x}_k)) + \nabla f(\mathbf{x}_k)\| \\ &\leq D \text{dist}(N_C(\mathbf{x}_k), -\nabla f(\mathbf{x}_k)), \end{aligned} \tag{39}$$

where in the first inequality we used  $(\mathbf{s}_k - \mathbf{x}_k)^\top [\pi_{N_C(\mathbf{x}_k)}(-\nabla f(\mathbf{x}_k))] \leq 0$  together with the Cauchy-Schwarz inequality, and  $\|\mathbf{s}_k - \mathbf{x}_k\| \leq D$  in the second.

### 5.2 $\mathcal{O}(1/k)$ rate for convex objectives

If  $f$  is non-convex, it is possible to prove a  $\mathcal{O}(1/\sqrt{k})$  rate for  $\min_{i \in [1:k]} G(x_i)$  (see, e.g., Lacoste-Julien 2016). On the other hand, if  $f$  is convex, we have an  $\mathcal{O}(1/k)$  rate on the optimality gap (see, e.g., Frank and Wolfe 1956; Levitin and Polyak 1966) for all the stepsizes discussed in Sect 4. Here we include a proof for the Lipschitz constant dependent stepsize  $\alpha_k$  given by (28).

**Theorem 1** *If  $f$  is convex and the stepsize is given by (28), then for every  $k \geq 1$*

$$f(\mathbf{x}_k) - f^* \leq \frac{2LD^2}{k + 2}. \tag{40}$$

Before proving the theorem we prove a lemma concerning the decrease of the objective in the case of a full FW step, that is a step with  $\mathbf{d}_k = \mathbf{d}_k^{FW}$  and with  $\alpha_k$  equal to 1, the maximal feasible stepsize.

**Lemma 2** *If  $\alpha_k = 1$  and  $\mathbf{d}_k = \mathbf{d}_k^{FW}$  then*

$$f(\mathbf{x}_{k+1}) - f^* \leq \frac{1}{2} \min \left\{ L \|\mathbf{d}_k\|^2, f(\mathbf{x}_k) - f^* \right\}. \tag{41}$$

**Proof** If  $\alpha_k = 1 = \alpha_k^{\max}$  then by Definitions (3) and (32)

$$G(\mathbf{x}_k) = -\nabla f(\mathbf{x}_k)^\top \mathbf{d}_k \geq L \|\mathbf{d}_k\|^2, \tag{42}$$

the last inequality following by Definition (28) and the assumption that  $\alpha_k = 1$ . By the standard Descent Lemma it also follows

$$f(\mathbf{x}_{k+1}) - f^* = f(\mathbf{x}_k + \mathbf{d}_k) - f^* \leq f(\mathbf{x}_k) - f^* + \nabla f(\mathbf{x}_k)^\top \mathbf{d}_k + \frac{L}{2} \|\mathbf{d}_k\|^2. \tag{43}$$

Considering the definition of  $\mathbf{d}_k$  and convexity of  $f$ , we get

$$f(\mathbf{x}_k) - f^* + \nabla f(\mathbf{x}_k)^\top \mathbf{d}_k \leq f(\mathbf{x}_k) - f^* + \nabla f(\mathbf{x}_k)^\top (\mathbf{x}^* - \mathbf{x}_k) \leq 0,$$

so that (43) entails  $f(\mathbf{x}_{k+1}) - f^* \leq \frac{L}{2} \|\mathbf{d}_k\|^2$ . To conclude, it suffices to apply to the RHS of (43) the inequality

$$f(\mathbf{x}_k) - f^* + \nabla f(\mathbf{x}_k)^\top \mathbf{d}_k + \frac{L}{2} \|\mathbf{d}_k\|^2 \leq f(\mathbf{x}_k) - f^* - \frac{1}{2} G(\mathbf{x}_k) \leq \frac{f(\mathbf{x}_k) - f^*}{2} \tag{44}$$

where we used (42) in the first inequality and  $G(\mathbf{x}_k) \geq f(\mathbf{x}_k) - f^*$  in the second.  $\square$

We can now proceed with the proof of the main result.

**Proof of Theorem 1** For  $k = 0$  and  $\alpha_0 = 1$  then by Lemma 2

$$f(\mathbf{x}_1) - f^* \leq \frac{L \|\mathbf{d}_0\|^2}{2} \leq \frac{LD^2}{2}. \tag{45}$$

If  $\alpha_0 < 1$  then

$$f(\mathbf{x}_0) - f^* \leq G(\mathbf{x}_0) < L \|\mathbf{d}_0\|^2 \leq LD^2. \tag{46}$$

Therefore in both cases (30) holds for  $k = 0$ .

Reasoning by induction, if (40) holds for  $k$  with  $\alpha_k = 1$ , then the claim is clear by (41).

On the other hand, if  $\alpha_k < \alpha_k^{\max} = 1$  then by Lemma 1, we have

$$\begin{aligned} f(\mathbf{x}_{k+1}) - f^* &\leq f(\mathbf{x}_k) - f^* - \frac{1}{2L} (\nabla f(\mathbf{x}_k)^\top \widehat{\mathbf{d}}_k)^2 \\ &\leq f(\mathbf{x}_k) - f^* - \frac{(\nabla f(\mathbf{x}_k)^\top \mathbf{d}_k)^2}{2LD^2} \\ &\leq f(\mathbf{x}_k) - f^* - \frac{(f(\mathbf{x}_k) - f^*)^2}{2LD^2} \\ &= (f(\mathbf{x}_k) - f^*) \left( 1 - \frac{f(\mathbf{x}_k) - f^*}{2LD^2} \right) \leq \frac{2LD^2}{k+3}, \end{aligned} \tag{47}$$

where we used  $\|d_k\| \leq D$  in the second inequality,  $\nabla f(x_k)^\top d_k = G(x_k) \geq f(x_k) - f^*$  in the third inequality; and the last inequality follows by induction hypothesis.  $\square$

As can be easily seen from above argument, the convergence rate of  $\mathcal{O}(1/k)$  is true also in more abstract normed spaces than  $\mathbb{R}^n$ , e.g. when  $C$  is a convex and weakly compact subset of a Banach space (see, e.g., Demyanov and Rubinov 1970; Dunn and Harshbarger 1978). A generalization for some unbounded sets is given in Ferreira and Sosa (2021). The bound is tight due to a zigzagging behaviour of the method near solutions on the boundary, leading to a rate of  $\Omega(1/k^{1+\delta})$  for every  $\delta > 0$  (see Canon and Cullum 1968 for further details), when the objective is a strictly convex quadratic function and the domain is a polytope.

Also the minimum FW gap  $\min_{i \in [0:k]} G(x_i)$  converges at a rate of  $\mathcal{O}(1/k)$  (see Jaggi 2013; Freund and Grigas 2016). In Freund and Grigas (2016), a broad class of step-sizes is examined, including  $\alpha_k = \frac{1}{k+1}$  and  $\alpha_k = \bar{\alpha}$  constant. For these step-sizes a convergence rate of  $\mathcal{O}\left(\frac{\ln(k)}{k}\right)$  is proved.

### 5.3 Variants

Active set FW variants mostly aim to improve over the  $\mathcal{O}(1/k)$  rate and also ensure support identification in finite time. They generate a sequence of active sets  $\{A_k\}$ , such that  $x_k \in \text{conv}(A_k)$ , and define alternative directions making use of these active sets.

For the *pairwise FW (PFW)* and the *away step FW (AFW)* (see Clarkson 2010; Lacoste-Julien and Jaggi 2015) we have that  $A_k$  must always be a subset of  $S_k$ , with  $x_k$  a convex combination of the elements in  $A_k$ . The away vertex  $v_k$  is then defined by

$$v_k \in \arg \max_{y \in A_k} \nabla f(x_k)^\top y. \tag{48}$$

The AFW direction, introduced in Wolfe (1970), is hence given by

$$\begin{aligned} d_k^{AS} &= x_k - v_k \\ d_k &\in \arg \max \{-\nabla f(x_k)^\top d : d \in \{d_k^{AS}, d_k^{FW}\}\}, \end{aligned} \tag{49}$$

while the PFW direction, as defined in Lacoste-Julien and Jaggi (2015) and inspired by the early work (Mitchell et al. 1974), is

$$d_k^{PFW} = d_k^{FW} + d_k^{AS} = s_k - v_k, \tag{50}$$

with  $s_k$  defined in (3).

The *FW method with in-face directions (FDW)* (see Freund et al. 2017; Guélat and Marcotte 1986), also known as Decomposition invariant Conditional Gradient (DiCG) when applied to polytopes (Bashiri and Zhang 2017), is defined exactly as the AFW, but with the minimal face  $\mathcal{F}(x_k)$  of  $C$  containing  $x_k$  as the active set. The *extended FW (EFW)* was introduced in Holloway (1974) and is also known as simplicial decomposition (Von Hohenbalken 1977). At every iteration the method

**Table 1** FW method and variants covered in this review

Variant	Direction	Active set
FW	$d_k = d_k^{FW} = s_k - x_k, \quad s_k \in \arg \max\{\nabla f(x_k)^\top x : x \in C\}$	–
AFW	$d_k \in \arg \max\{-\nabla f(x_k)^\top d : d \in \{x_k - v_k, d_k^{FW}\}, v_k \in A_k\}$	$A_{k+1} \subseteq A_k \cup \{s_k\}$
PFW	$d_k = s_k - v_k, \quad v_k \in \arg \max\{\nabla f(x_k)^\top v_k : v_k \in A_k\}$	$A_{k+1} \subseteq A_k \cup \{s_k\}$
EFW	$d_k = y_k - x_k, \quad y_k \in \arg \min\{f(y) : y \in \text{conv}(A_k)\}$	$A_{k+1} \subseteq A_k \cup \{s_k\}$
DFW	$d_k \in \arg \max\{-\nabla f(x_k)^\top d : d \in \{x_k - v_k, d_k^{FW}\}, v_k \in A_k\}$	$A_k = \mathcal{F}(x_k)$

minimizes the objective in the current active set  $A_{k+1}$

$$x_{k+1} \in \arg \min_{y \in \text{conv}(A_{k+1})} f(y), \tag{51}$$

where  $A_{k+1} \subseteq A_k \cup \{s_k\}$  (see, e.g., Clarkson 2010, Algorithm 4.2). A more general version of the EFW, only approximately minimizing on the current active set, was introduced in Lacoste-Julien and Jaggi (2015) under the name of fully corrective FW. In Table 1, we report the main features of the classic FW and of the variants under analysis.

### 5.4 Sparse approximation properties

As discussed in the previous section, for the classic FW method and the AFW, PFW, EFW variants  $x_k$  can always be written as a convex combination of elements in  $A_k \subset S_k$ , with  $|A_k| \leq k + 1$ . Even for the DFW we still have the weaker property that  $x_k$  must be an affine combination of elements in  $A_k \subset A$  with  $|A_k| \leq k + 1$ . It turns out that the convergence rate of methods with this property is  $\Omega(\frac{1}{k})$  in high dimension. More precisely, if  $C = \text{conv}(A)$  with  $A$  compact, the  $\mathcal{O}(1/k)$  rate of the classic FW method is worst case optimal given the sparsity constraint

$$x_k \in \text{aff}(A_k) \text{ with } A_k \subset A, \quad |A_k| \leq k + 1. \tag{52}$$

An example where the  $\mathcal{O}(1/k)$  rate is tight was presented in Jaggi (2013). Let  $C = \Delta_{n-1}$  and  $f(x) = \|x - \frac{1}{n} e\|^2$ . Clearly,  $f^* = 0$  with  $x^* = \frac{1}{n} e$ . Then it is easy to see that  $\min\{f(x) - f^* : \|x\|_0 \leq k + 1\} \geq \frac{1}{k+1} - \frac{1}{n}$  for every  $k \in \mathbb{N}$ , so that in particular under (52) with  $A_k = \{e_i : i \in [1:n]\}$ , the rate of any FW variant must be  $\Omega(\frac{1}{k})$ .

### 5.5 Affine invariance

The FW method and the AFW, PFW, EFW are affine invariant (Jaggi 2013). More precisely, let  $P$  be a linear transformation,  $\hat{f}$  be such that  $\hat{f}(Px) = f(x)$  and  $\hat{C} = P(C)$ . Then for every sequence  $\{x_k\}$  generated by the methods applied to  $(f, C)$ , the sequence  $\{y_k\} := \{Px_k\}$  can be generated by the FW method with the same stepsizes applied to  $(\hat{f}, \hat{C})$ . As a corollary, considering the special case where  $P$  is the matrix collecting

the elements of  $A$  as columns, one can prove results on  $C = \Delta_{|A|-1}$  and generalize them to  $\hat{C} := \text{conv}(A)$  by affine invariance.

An affine invariant convergence rate bound for convex objectives can be given using the curvature constant

$$\kappa_{f,C} := \sup \left\{ 2 \frac{f(\alpha y + (1-\alpha)x) - f(x) - \alpha \nabla f(x)^\top (y-x)}{\alpha^2} : \{x, y\} \subset C, \alpha \in (0, 1] \right\}. \tag{53}$$

It is easy to prove that  $\kappa_{f,C} \leq LD^2$  if  $D$  is the diameter of  $C$ . In the special case where  $C = \Delta_{n-1}$  and  $f(x) = x^\top \tilde{A}^\top \tilde{A} x + b^\top x$ , then  $\kappa_{f,C} \leq \text{diam}(A \Delta_{n-1})^2$  for  $A^\top = [\tilde{A}^\top, b]$ ; see Clarkson (2010).

When the method uses the stepsize sequence (25), it is possible to give the following affine invariant convergence rate bounds (see Freund and Grigas 2016):

$$\begin{aligned} f(x_k) - f^* &\leq \frac{2\kappa_{f,C}}{k+4}, \\ \min_{i \in [1:k]} G(x_i) &\leq \frac{9\kappa_{f,C}}{2k}, \end{aligned} \tag{54}$$

thus in particular slightly improving the rate we gave in Theorem 1 since we have that  $\kappa_{f,C} \leq LD^2$ .

### 5.6 Support identification for the AFW

It is a classic result that the AFW under some strict complementarity conditions and for strongly convex objectives identifies in finite time the face containing the solution (Guélat and Marcotte 1986). Here we report some explicit bounds for this property proved in Bomze et al. (2020). We first assume that  $C = \Delta_{n-1}$ , and introduce the multiplier functions

$$\lambda_i(x) = \nabla f(x)^\top (e_i - x) \tag{55}$$

for  $i \in [1:n]$ . Let  $x^*$  be a stationary point for  $f$ , with the objective  $f$  not necessarily convex. It is easy to check that  $\{\lambda_i(x^*)\}_{i \in [1:n]}$  coincide with the Lagrangian multipliers. Furthermore, by complementarity conditions we have  $x_i^* \lambda_i(x^*) = 0$  for every  $i \in [1:n]$ . It follows that the set

$$I(x^*) := \{i \in [1:n] : \lambda_i(x^*) = 0\}$$

contains the support of  $x^*$ ,

$$\text{supp}(x^*) := \{i \in [1:n] : x_i^* > 0\}.$$

The next lemma uses  $\lambda_i$ , and the Lipschitz constant  $L$  of  $\nabla f$ , to give a lower bound of the so-called *active set radius*  $r_*$ , defining a neighborhood of  $x^*$ . Starting the algorithm

in this neighbourhood, the active set (the minimal face of  $C$  containing  $x^*$ ) is identified in a limited number of iterations.

**Lemma 3** *Let  $x^*$  be a stationary point for  $f$  on the boundary of  $\Delta_{n-1}$ ,  $\delta_{\min} = \min_{i:\lambda_i(x^*)>0} \lambda_i(x^*)$  and*

$$r_* = \frac{\delta_{\min}}{\delta_{\min} + 2L}. \tag{56}$$

*Assume that for every  $k$  for which  $d_k = d_k^A$  holds, the step size  $\alpha_k$  is not smaller than the stepsize given by (28),  $\alpha_k(L) \leq \alpha_k$ .*

*If  $\|x_k - x^*\|_1 < r_*$ , then for some*

$$j \leq \min\{n - |I(x^*)|, |\text{supp}(x_k)| - 1\}$$

*we have  $\text{supp}(x_{k+j}) \subseteq I(x^*)$  and  $\|x_{k+j} - x^*\|_1 < r_*$ .*

**Proof** Follows from (Bomze et al. 2020, Theorem 3.3), since under the assumptions the AFW sets one variable in  $\text{supp}(x_k) \setminus I(x^*)$  to zero at every step without increasing the 1-norm distance from  $x^*$ . □

The above lemma does not require convexity and was applied in Bomze et al. (2020) to derive active set identification bounds in several convex and non-convex settings. Here we focus on the case where the domain  $C = \text{conv}(A)$  with  $|A| < +\infty$  is a generic polytope, and where  $f$  is  $\mu$ -strongly convex for some  $\mu > 0$ , i.e.

$$f(y) \geq f(x) + \nabla f(x)^\top(y - x) + \frac{\mu}{2}\|x - y\|^2 \quad \text{for all } \{x, y\} \subset C. \tag{57}$$

Let  $E_C(x^*)$  be the face of  $C$  exposed by  $\nabla f(x^*)$ :

$$E_C(x^*) := \arg \min_{x \in C} \nabla f(x^*)^\top x, \tag{58}$$

Let then  $\theta_A$  be the Hoffman constant (see Beck and Shtern 2017) related to  $[\bar{A}^\top, I_n, e, -e]^\top$ , with  $\bar{A}$  the matrix having as columns the elements in  $A$ . Finally, consider the function  $f_A(y) := f(\bar{A}y)$  on  $\Delta_{|A|-1}$ , and let  $L_A$  be the Lipschitz constant of  $\nabla f_A$  as well as

$$\delta_{\min} := \min_{a \in A \setminus E_C(x^*)} \nabla f(x^*)^\top(a - x^*) \quad \text{and} \quad r_*(x^*) := \frac{\delta_{\min}}{\delta_{\min} + 2L_A}.$$

Using linearity of AFW convergence for strongly convex objectives (see Sect. 6.1), we have the following result:

**Theorem 2** *The sequence  $\{x_k\}$  generated by the AFW with  $x_0 \in A$  enters  $E_C(x^*)$  for*

$$k \geq \max \left\{ 2 \frac{\ln(h_0) - \ln(\mu_A r_*(x^*)^2/2)}{\ln(1/q)}, 0 \right\}, \tag{59}$$

where  $\mu_A = \frac{\mu}{n\theta_A^2}$  and  $q \in (0, 1)$  is the constant related to the linear convergence rate of the AFW, i.e.  $h_k \leq q^k h_0$  for all  $k$ .

**Proof of Theorem 2** (sketch) We present an argument in the case  $C = \Delta_{n-1}$ ,  $A = \{e_i\}_{i \in [1:n]}$  which can be easily extended by affine invariance to the general case (see Bomze et al. 2020 for details). In this case  $\theta_A \geq 1$  and we can define  $\bar{\mu} := \mu/n \geq \mu_A$ . To start with, the number of steps needed to reach the condition

$$h_k \leq \frac{\mu}{2n} r_*(x^*)^2 = \frac{\bar{\mu}}{2} r_*(x^*)^2 \tag{60}$$

is at most

$$\bar{k} = \max \left\{ \left\lceil \frac{\ln(h_0) - \ln(\bar{\mu} r_*(x^*)^2/2)}{\ln(1/q)} \right\rceil, 0 \right\}.$$

Now we combine  $n\|\cdot\| \geq \|\cdot\|_1$  with strong convexity and relation (60) to obtain  $\|x_k - x^*\|_1 \leq r_*(x^*)$ , hence in particular  $\|x_k - x^*\|_1 \leq r_*(x^*)$  for every  $k \geq \bar{k}$ . Since  $x_0$  is a vertex of the simplex, and at every step at most one coordinate is added to the support of the current iterate,  $|\text{supp}(x_{\bar{k}})| \leq \bar{k} + 1$ . The claim follows by applying Lemma 3.  $\square$

Additional bounds under a quadratic growth condition weaker than strong convexity and strict complementarity are reported in Garber (2020).

Convergence and finite time identification for the PFW and the AFW are proved in Bomze et al. (2019) for a specific class of non-convex minimization problems over the standard simplex, under the additional assumption that the sequence generated has a finite set of limit points. In another line of work, active set identification strategies combined with FW variants have been proposed in Cristofari et al. (2020) and Sun (2020).

### 5.7 Inexact linear oracle

In many real-world applications, linear subproblems can only be solved approximately. This is the reason why the convergence of FW variants is often analyzed under some error term for the linear minimization oracle (see, e.g., Braun et al. 2019, 2017; Freund and Grigas 2016; Jaggi 2013; Konnov 2018). A common assumption, relaxing the FW vertex exact minimization property, is to have access to a point (usually a vertex)  $\tilde{s}_k$  such that

$$\nabla f(x_k)^\top (\tilde{s}_k - x_k) \leq \min_{s \in C} \nabla f(x_k)^\top (s - x_k) + \delta_k, \tag{61}$$

for a sequence  $\{\delta_k\}$  of non negative approximation errors.

If the sequence  $\{\delta_k\}$  is constant and equal to some  $\delta > 0$ , then trivially the lowest possible approximation error achieved by the FW method is  $\delta$ . At the same time,

(Freund and Grigas 2016, Theorem 5.1) implies a rate of  $\mathcal{O}(\frac{1}{k} + \delta)$  if the stepsize  $\alpha_k = \frac{2}{k+2}$  is used.

The  $\mathcal{O}(1/k)$  rate can be instead retrieved by assuming that  $\{\delta_k\}$  converges to 0 quickly enough, and in particular if

$$\delta_k = \frac{\delta \kappa_{f,C}}{k + 2} \tag{62}$$

for a constant  $\delta > 0$ . Under (62), in Jaggi (2013) a convergence rate of

$$f(x_k) - f^* \leq \frac{2\kappa_{f,C}}{k + 2} (1 + \delta) \tag{63}$$

was proved for the FW method with  $\alpha_k$  given by exact line search or equal to  $\frac{2}{k+2}$ , as well as for the EFW.

A variant making use of an approximated linear oracle recycling previous solutions to the linear minimization subproblem is studied in Braun et al. (2019). In Freund and Grigas (2016), Hogan (1971), the analysis of the classic FW method is extended to the case of inexact gradient information. In particular in Freund and Grigas (2016), assuming the availability of the  $(\delta, L)$  oracle introduced in Devolder et al. (2014), a convergence rate of  $\mathcal{O}(1/k + \delta k)$  is proved.

## 6 Improved rates for strongly convex objectives

### 6.1 Linear convergence under an angle condition

In the rest of this section we assume that  $f$  is  $\mu$ -strongly convex (57). We also assume that the stepsize is given by exact linesearch or by (28).

Under the strong convexity assumption, an asymptotic linear convergence rate for the FDFW on polytopes was given in the early work (Guélat and Marcotte 1986). Furthermore, in Garber and Hazan (2016) a linearly convergent variant was proposed, making use however of an additional local linear minimization oracle. See also Table 2 for a list of improvements on the  $\mathcal{O}(1/k)$  rate under strong convexity.

**Table 2** Known convergence rates for the FW method and the variants covered in this review

Method	Objective	Domain	Assumptions	Rate
FW	NC	Generic	–	$\mathcal{O}(1/\sqrt{k})$
FW	C	Generic	–	$\mathcal{O}(1/k)$
FW	SC	Generic	$x^* \in \text{ri}(C)$	Linear
Variants	SC	Polytope	–	Linear
FW	SC	Strongly convex	–	$\mathcal{O}(1/k^2)$
FW	SC	Strongly convex	$\min \ \nabla f(x)\  > 0$	Linear

NC, C and SC stand for non-convex, convex and strongly convex respectively

Recent works obtain linear convergence rates by proving the angle condition

$$-\nabla f(x_k)^\top \widehat{d}_k \geq \frac{\tau}{\|x_k - x^*\|} \nabla f(x_k)^\top (x_k - x^*) \tag{64}$$

for some  $\tau > 0$  and some  $x^* \in \arg \min_{x \in C} f(x)$ . As we shall see in the next lemma, under (64) it is not difficult to prove linear convergence rates in the number of *good steps*. These are FW steps with  $\alpha_k = 1$  and steps in any descent direction with  $\alpha_k < 1$ .

**Lemma 4** *If the step  $k$  is a good step and (64) holds, then*

$$h_{k+1} \leq \max \left\{ \frac{1}{2}, 1 - \frac{\tau^2 \mu}{L} \right\} h_k. \tag{65}$$

**Proof** If the step  $k$  is a full FW step then Lemma 2 entails  $h_{k+1} \leq \frac{1}{2} h_k$ . In the remaining case, first observe that by strong convexity

$$\begin{aligned} f^* = f(x^*) &\geq f(x_k) + \nabla f(x_k)^\top (x^* - x_k) + \frac{\mu}{2} \|x_k - x^*\|^2 \\ &\geq \min_{\alpha \in \mathbb{R}} \left[ f(x_k) + \alpha \nabla f(x_k)^\top (x^* - x_k) + \frac{\alpha^2 \mu}{2} \|x_k - x^*\|^2 \right] \\ &= f(x_k) - \frac{1}{2\mu \|x_k - x^*\|^2} \left[ \nabla f(x_k)^\top (x_k - x^*) \right]^2, \end{aligned} \tag{66}$$

which means

$$h_k \leq \frac{1}{2\mu \|x_k - x^*\|^2} \left[ \nabla f(x_k)^\top (x_k - x^*) \right]^2. \tag{67}$$

We can then proceed using the bound (30) from Lemma 1 in the following way:

$$\begin{aligned} h_{k+1} = f(x_{k+1}) - f^* &\leq f(x_k) - f^* - \frac{1}{2L} \left[ \nabla f(x_k)^\top \widehat{d}_k \right]^2 \\ &\leq h_k - \frac{\tau^2}{2L \|x_k - x^*\|^2} \left[ \nabla f(x_k)^\top (x_k - x^*) \right]^2 \\ &\leq h_k \left( 1 - \frac{\tau^2 \mu}{L} \right), \end{aligned} \tag{68}$$

where we used (64) in the second inequality and (67) in the third one. □

As a corollary, under (64) we have the rate

$$f(x_k) - f^* = h_k \leq \max \left\{ \frac{1}{2}, 1 - \frac{\tau^2 \mu}{L} \right\}^{\gamma(k)} h_0 \tag{69}$$

for any method with non increasing  $\{f(x_k)\}$  and following Algorithm 1, with  $\gamma(k) \leq k$  an integer denoting the number of good steps until step  $k$ . It turns out that for all the variants we introduced in this review we have  $\gamma(k) \geq Tk$  for some constant  $T > 0$ . When  $x^*$  is in the relative interior of  $C$ , the FW method satisfies (64) and we have the following result (see Guélat and Marcotte 1986; Lacoste-Julien and Jaggi 2015):

**Theorem 3** *If  $x^* \in \text{ri}(C)$ , then*

$$f(x_k) - f^* \leq \left[ 1 - \frac{\mu}{L} \left( \frac{\text{dist}(x^*, \partial C)}{D} \right)^2 \right]^k (f(x_0) - f^*). \tag{70}$$

**Proof** We can assume for simplicity  $\text{int}(C) \neq \emptyset$ , since otherwise we can restrict ourselves to the affine hull of  $C$ . Let  $\delta = \text{dist}(x^*, \partial C)$  and  $g = -\nabla f(x_k)$ . First, by assumption we have  $x^* + \delta \widehat{g} \in C$ . Therefore

$$g^T d_k^{FW} \geq g^T ((x^* + \delta \widehat{g}) - x) = \delta g^T \widehat{g} + g^T (x^* - x) \geq \delta \|g\| + f(x) - f^* \geq \delta \|g\|, \tag{71}$$

where we used  $x^* + \delta \widehat{g} \in C$  in the first inequality and convexity in the second. We can conclude

$$g^T \frac{d_k^{FW}}{\|d_k^{FW}\|} \geq g^T \frac{d_k^{FW}}{D} \geq \frac{\delta}{D} \|g\| \geq \frac{\delta}{D} g^T \left( \frac{x_k - x^*}{\|x_k - x^*\|} \right). \tag{72}$$

The thesis follows by Lemma 4, noticing that for  $\tau = \frac{\text{dist}(x^*, \partial C)}{D} \leq \frac{1}{2}$  we have  $1 - \tau^2 \frac{\mu}{L} > \frac{1}{2}$ . □

In Lacoste-Julien and Jaggi (2015), the authors proved that directions generated by the AFW and the PFW on polytopes satisfy condition (64), with  $\tau = \text{PWidth}(A)/D$  and  $\text{PWidth}(A)$ , pyramidal width of  $A$ . While  $\text{PWidth}(A)$  was originally defined with a rather complex minmax expression, in Peña and Rodriguez (2018) it was then proved

$$\text{PWidth}(A) = \min_{F \in \text{faces}(C)} \text{dist}(F, \text{conv}(A \setminus F)). \tag{73}$$

This quantity can be explicitly computed in a few special cases. For  $A = \{0, 1\}^n$  we have  $\text{PWidth}(A) = 1/\sqrt{n}$ , while for  $A = \{e_i\}_{i \in [1:n]}$  (so that  $C$  is the  $n - 1$  dimensional simplex)

$$\text{PWidth}(A) = \begin{cases} \frac{2}{\sqrt{n}} & \text{if } n \text{ is even} \\ \frac{2}{\sqrt{n-1/n}} & \text{if } n \text{ is odd.} \end{cases} \tag{74}$$

Angle conditions like (64) with  $\tau$  dependent on the number of vertices used to represent  $x_k$  as a convex combination were given in Bashiri and Zhang (2017) and Beck and Shtern (2017) for the FDFW and the PFW respectively. In particular, in Beck and Shtern (2017) a geometric constant  $\Omega_C$  called vertex-facet distance was defined as

$$\Omega_C = \min\{\text{dist}(v, H) : v \in V(C), H \in \mathcal{H}(C), v \notin H\}, \tag{75}$$

with  $V(C)$  the set of vertices of  $C$ , and  $\mathcal{H}(C)$  the set of supporting hyperplanes of  $C$  (containing a facet of  $C$ ). Then condition (64) is satisfied for  $\tau = \Omega_C/s$ , with  $d_k$  the PFW direction and  $s$  the number of points used in the active set  $A_k$ .

In Bashiri and Zhang (2017), a geometric constant  $H_s$  was defined depending on the minimum number  $s$  of vertices needed to represent the current point  $x_k$ , as well as on the proper<sup>2</sup> inequalities  $q_i^\top x \leq b_i$ ,  $i \in [1:m]$ , appearing in a description of  $C$ . For each of these inequalities the *second gap*  $g_i$  was defined as

$$g_i = \max_{v \in V(C)} q_i^\top v - \text{secondmax}_{v \in V(C)} q_i^\top v, \quad i \in [1:m], \quad (76)$$

with the *secondmax* function giving the second largest value achieved by the argument. Then  $H_s$  is defined as

$$H_s := \max \left\{ \sum_{j=1}^n \left( \sum_{i \in S} \frac{a_{ij}}{g_i} \right)^2 : S \in \binom{[1:m]}{s} \right\}. \quad (77)$$

The arguments used in the paper imply that (64) holds with  $\tau = \frac{1}{2D\sqrt{H_s}}$  if  $d_k$  is a FDFW direction and  $x_k$  the convex combination of at most  $s$  vertices. We refer the reader to Peña and Rodríguez (2018) and Rademacher and Shu (2020) for additional results on these and related constants.

The linear convergence results for strongly convex objectives are extended to compositions of strongly convex objectives with affine transformations in Beck and Shtern (2017), Lacoste-Julien and Jaggi (2015), Peña and Rodríguez (2018). In Gutman and Pena (2021), the linear convergence results for the AFW and the FW method with minimum in the interior are extended with respect to a generalized condition number  $L_{f,C,D}/\mu_{f,C,D}$ , with  $D$  a distance function on  $C$ .

For the AFW, the PFW and the FDFW, linear rates with no bad steps ( $\gamma(k) = k$ ) are given in Rinaldi and Zeffiro (2020) for non-convex objectives satisfying a Kurdyka-Łojasiewicz inequality. In Rinaldi and Zeffiro (2020), condition (64) was proved for the FW direction and orthographic retractions on some convex sets with smooth boundary. The work Combettes and Pokutta (2020) introduces a new FW variant using a subroutine to align the descent direction with the projection on the tangent cone of the negative gradient, thus implicitly maximizing  $\tau$  in (64).

## 6.2 Strongly convex domains

When  $C$  is strongly convex we have a  $\mathcal{O}(1/k^2)$  rate (see, e.g., Garber and Hazan 2015; Kerdreux et al. 2021) for the classic FW method. Furthermore, when  $C$  is  $\beta_C$ -strongly convex and  $\|\nabla f(x)\| \geq c > 0$ , then we have the linear convergence rate (see Demyanov and Rubinov 1970; Dunn 1979; Kerdreux et al. 2020; Levitin and Polyak 1966)

$$h_{k+1} \leq \max \left\{ \frac{1}{2}, 1 - \frac{L}{2c\beta_C} \right\} h_k. \quad (78)$$

Finally, it is possible to interpolate between the  $\mathcal{O}(1/k^2)$  rate of the strongly convex setting and the  $\mathcal{O}(1/k)$  rate of the general convex one by relaxing strong convexity of

<sup>2</sup> i.e., those inequalities strictly satisfied for some  $x \in C$ .

the objective with Hölderian error bounds (Xu and Yang 2018) and also by relaxing strong convexity of the domain with uniform convexity (Kerdreux et al. 2021).

## 7 Extensions

### 7.1 Block coordinate Frank–Wolfe method

The block coordinate FW (BCFW) was introduced in Lacoste-Julien et al. (2013) for block product domains of the form  $C = C^{(1)} \times \dots \times C^{(m)} \subseteq \mathbb{R}^{n_1 + \dots + n_m}$ , and applied to structured SVM training. The algorithm operates by selecting a random block and performing a FW step in that block. Formally, for  $s \in \mathbb{R}^{m_i}$  let  $s^{(i)} \in \mathbb{R}^n$  be the vector with all blocks equal to  $\mathbf{o}$  except for the  $i$ -th block equal to  $s$ . We can write the direction of the BCFW as

$$\begin{aligned} \mathbf{d}_k &= s_k^{(i)} - \mathbf{x}_k \\ s_k &\in \arg \min_{s \in C^{(i)}} \nabla f(\mathbf{x}_k)^\top s \end{aligned} \tag{79}$$

for a random index  $i \in [1 : n]$ .

In Lacoste-Julien et al. (2013), a convergence rate of

$$\mathbb{E}[f(x_k)] - f^* \leq \frac{2Km}{k + 2m} \tag{80}$$

is proved, for  $K = h_0 + \kappa_f^\otimes$ , with  $\kappa_f^\otimes$  the product domain curvature constant, defined as  $\kappa_f^\otimes = \sum \kappa_f^{\otimes, i}$  where  $\kappa_f^{\otimes, i}$  are the curvature constants of the objective fixing the blocks outside  $C^{(i)}$ :

$$\kappa_f^{\otimes, i} := \sup \left\{ 2 \frac{f(\mathbf{x} + \alpha \mathbf{d}^{(i)}) - f(\mathbf{x}) - \alpha \nabla f(\mathbf{x})^\top \mathbf{d}^{(i)}}{\alpha^2} : \mathbf{d} \in C - \mathbf{x}, \mathbf{x} \in C, \alpha \in (0, 1] \right\}. \tag{81}$$

An asynchronous and parallel generalization for this method was given in Wang et al. (2016). This version assumes that a cloud oracle is available, modeling a set of worker nodes each sending information to a server at different times. This information consists of an index  $i$  and the following LMO on  $C^{(i)}$ :

$$s_{(i)} \in \arg \min_{s \in C^{(i)}} \nabla f(\mathbf{x}_{\tilde{k}})^\top s^{(i)}. \tag{82}$$

The algorithm is called asynchronous because  $\tilde{k}$  can be smaller than  $k$ , modeling a delay in the information sent by the node. Once the server has collected a minibatch  $S$  of  $\tau$  distinct indexes (overwriting repetitions), the descent direction is defined as

$$\mathbf{d}_k = \sum_{i \in S} s_{(i)}^{(i)}, \tag{83}$$

If the indices sent by the nodes are i.i.d., then under suitable assumptions on the delay, a convergence rate of

$$\mathbb{E}[f(\mathbf{x}_k)] - f^* \leq \frac{2mK_\tau}{\tau^2k + 2m} \quad (84)$$

can be proved, where  $K_\tau = m\kappa_{f,\tau}^\otimes(1 + \delta) + h_0$  for  $\delta$  depending on the delay error, with  $\kappa_{f,\tau}^\otimes$  the average curvature constant in a minibatch keeping all the components not in the minibatch fixed.

In Osokin et al. (2016), several improvements are proposed for the BCFW, including an adaptive criterion to prioritize blocks based on their FW gap, and block coordinate versions of the AFW and the PFW variants.

In Shah et al. (2015), a multi plane BCFW approach is proposed in the specific case of the structured SVM, based on caching supporting planes in the primal, corresponding to block linear minimizers in the dual. In Berrada et al. (2018), the duality for structured SVM between BCFW and stochastic subgradient descent is exploited to define a learning rate schedule for neural networks based on only one hyper parameter. The block coordinate approach is extended to the generalized FW in Beck et al. (2015), with coordinates however picked in a cyclic order.

## 7.2 Variants for the min-norm point problem

Consider the min-norm point (MNP) problem

$$\min_{\mathbf{x} \in C} \|\mathbf{x}\|_* , \quad (85)$$

with  $C$  a closed convex subset of  $\mathbb{R}^n$  and  $\|\cdot\|_*$  a norm on  $\mathbb{R}^n$ . In Wolfe (1976), a FW variant is introduced to solve the problem when  $C$  is a polytope and  $\|\cdot\|_*$  is the standard Euclidean norm  $\|\cdot\|$ . Similarly to the variants introduced in Sect. 5.3, it generates a sequence of active sets  $\{A_k\}$  with  $s_k \in A_{k+1}$ . At the step  $k$  the norm is minimized on the affine hull  $\text{aff}(A_k)$  of the current active set  $A_k$ , that is

$$\mathbf{v}_k = \arg \min_{\mathbf{y} \in \text{aff}(A_k)} \|\mathbf{y}\| . \quad (86)$$

The descent direction  $\mathbf{d}_k$  is then defined as

$$\mathbf{d}_k = \mathbf{v}_k - \mathbf{x}_k , \quad (87)$$

and the stepsize is given by a tailored linesearch that allows to remove some of the atoms in the set  $A_k$  (see, e.g. Lacoste-Julien and Jaggi 2015; Wolfe 1976). Whenever  $\mathbf{x}_{k+1}$  is in the relative interior of  $\text{conv}(A_k)$ , the FW vertex is added to the active set (that is,  $s_k \in A_{k+1}$ ). Otherwise, at least one of the vertices not appearing in a convex representation of  $\mathbf{x}_k$  is removed. This scheme converges linearly when applied to generic smooth strongly convex objectives (see, e.g., Lacoste-Julien and Jaggi 2015).

In Harchaoui et al. (2015), a FW variant is proposed for minimum norm problems of the form

$$\min\{\|x\|_* : f(x) \leq 0, x \in K\} \tag{88}$$

with  $K$  a convex cone,  $f$  convex with  $L$ -Lipschitz gradient. In particular, the optimization domain is  $C = \{x \in \mathbb{R}^n : f(x) \leq 0\} \cap K$ . The technique proposed in the article applies the standard FW method to the problems

$$\min\{f(x) : \|x\|_* \leq \delta_k, x \in K\},$$

with  $\{\delta_k\}$  an increasing sequence convergent to the optimal value  $\bar{\delta}$  of the problem (88). Let  $C(\delta) = \{x \in \mathbb{R}^n : \|x\|_* \leq \delta\} \cap K$  for  $\delta \geq 0$ , and let

$$LM(r) \in \arg \min_{x \in C(1)} r^\top x,$$

so that by homogeneity for every  $k$  the linear minimization oracle on  $C(\delta_k)$  is given by

$$LMO_{C(\delta_k)}(r) = \delta_k LM(r). \tag{89}$$

For every  $k$ , applying the FW method with suitable stopping conditions an approximate minimizer  $x_k$  of  $f(x)$  over  $C(\delta_k)$  is generated, with an associated lower bound on the objective, an affine function in  $y$ :

$$f_k(y) := f(x_k) + \nabla f(x_k)^\top (y - x_k). \tag{90}$$

Then the function

$$\ell_k(\delta) := \min_{y \in C(\delta)} f_k(y) = f_k(\delta LM(g_k)) \quad \text{with } g_k = \nabla f(x_k) \tag{91}$$

is decreasing and affine in  $\delta$  and satisfies

$$\ell_k(\delta) = \min_{y \in C(\delta)} f_k(y) \leq F(\delta) := \min_{y \in C(\delta)} f(y). \tag{92}$$

Therefore, for

$$\bar{\ell}_k(\delta) = \max_{i \in [1:k]} \ell_i(\delta) \leq F(\delta)$$

the quantity  $\delta_{k+1}$  can be defined as  $\min\{\delta \geq 0 : \bar{\ell}_k(\delta) \leq 0\}$ , hence  $F(\delta_{k+1}) \geq 0$ . A complexity bound of  $\mathcal{O}(\frac{1}{\varepsilon} \ln(\frac{1}{\varepsilon}))$  was given to achieve precision  $\varepsilon$  applying this method, with  $\mathcal{O}(1/\varepsilon)$  iterations per subproblem and length of the sequence  $\{\delta_k\}$  at most  $\mathcal{O}(\ln(1/\varepsilon))$  (see (Harchaoui et al. 2015, Theorem 2) for details).

### 7.3 Variants for optimization over the trace norm ball

The FW method has found many applications for optimization problems over the trace norm ball. In this case, as explained in Example 3.4, linear optimization can be obtained by computing the top left and right singular vectors of the matrix  $-\nabla f(X_k)$ , an operation referred to as 1-SVD (see Allen-Zhu et al. 2017).

In the work Freund et al. (2017), the FDFW is applied to the matrix completion problem (13), thus generating a sequence of matrices  $\{X_k\}$  with  $\|X_k\|_* \leq \delta$  for every  $k$ . The method can be implemented efficiently exploiting the fact that for  $X$  on the boundary of the nuclear norm ball, there is a simple expression for the face  $\mathcal{F}(X)$ . For  $X \in \mathbb{R}^{m \times n}$  with  $\text{rank}(X) = k$  let  $UDV^T$  be the thin SVD of  $X$ , so that  $D \in \mathbb{R}^{k \times k}$  is the diagonal matrix of non zero singular values for  $X$ , with corresponding left and right singular vectors in the columns of  $U \in \mathbb{R}^{m \times k}$  and  $V \in \mathbb{R}^{n \times k}$  respectively. If  $\|X\|_* = \delta$  then the minimal face of the domain containing  $X$  is the set

$$\mathcal{F}(X) = \{X \in \mathbb{R}^{m \times n} : X = UMV^T \text{ for } M = M^T \text{ psd with } \|M\|_* = \delta\}. \quad (93)$$

It is not difficult to see that we have  $\text{rank}(X_k) \leq k + 1$  for every  $k \in \mathbb{N}$ , as well. Furthermore, the thin SVD of the current iterate  $X_k$  can be updated efficiently both after FW steps and after in face steps. The convergence rate of the FDFW in this setting is still  $\mathcal{O}(1/k)$ .

In the recent work Wang et al. (2020), an unbounded variant of the FW method is applied to solve a generalized version of the trace norm ball optimization problem:

$$\min_{X \in \mathbb{R}^{m \times n}} \{f(X) : \|PXQ\|_* \leq \delta\} \quad (94)$$

with  $P, Q$  singular matrices. The main idea of the method is to decompose the domain in the sum  $S + T$  between the kernel  $T$  of the linear function  $\varphi_{P,Q}(X) = PXQ$  and a bounded set  $S \subset T^\perp$ . Then gradient descent steps in the unbounded component  $T$  are alternated to FW steps in the bounded component  $S$ . The authors apply this approach to the generalized LASSO as well, using the AFW for the bounded component.

In Allen-Zhu et al. (2017), a variant of the classic FW using  $k$ -SVD (computing the top  $k$  left and right singular vectors for the SVD) is introduced, and it is proved that it converges linearly for strongly convex objectives when the solution has rank at most  $k$ . In Mu et al. (2016), the FW step is combined with a proximal gradient step for a quadratic problem on the product of the nuclear norm ball with the  $\ell_1$  ball. Approaches using an equivalent formulation on the spectrahedron introduced in Jaggi and Sulovský (2010) are analyzed in Ding et al. (2020), Garber (2019).

## 8 Conclusions

While the concept of the FW method is quite easy to understand, its advantages, witnessed by a multitude of related work, may not be apparent to someone not closely familiar with the subject. Therefore we considered, in Sect. 3, several motivating applications, ranging from classic optimization to more recent machine learning problems.

As in any line search-based method, the proper choice of stepsize is an important ingredient to achieve satisfactory performance. In Sect. 4, we review several options for stepsizes in first order methods, which are closely related both to the theoretical analysis as well as to practical implementation issues, guaranteeing fast convergence. This scope was investigated in more detail in Sect. 5 covering main results about the FW method and its most popular variants, including the  $\mathcal{O}(1/k)$  convergence rate for convex objectives, affine invariance, the sparse approximation property, and support identification. The account is complemented by a report on recent progress in improving on the  $\mathcal{O}(1/k)$  convergence rate in Sect. 6. Versatility and efficiency of this approach was also illustrated in the final Sect. 7 describing present recent FW variants fitting different optimization frameworks and computational environments, in particular block coordinate, distributed, parametrized, and trace norm optimization. For sure many other interesting and relevant aspects of FW and friends could not find their way into this review because of space and time limitations, but the authors hope to have convinced readers that FW merits a consideration even by non-experts in first-order optimization.

**Acknowledgements** The authors would like to thank an anonymous referee for their diligence and the Editors-in-Chief for their trust and encouragement, offering us the opportunity to publish this Invited Review.

**Funding** Open access funding provided by University of Vienna.

## Declarations

**Conflict of interest** The authors confirm there is no conflict of interest.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

- Ahipaşaoglu SD, Sun P, Todd MJ (2008) Linear convergence of a modified Frank–Wolfe algorithm for computing minimum-volume enclosing ellipsoids. *Optim Methods Soft* 23(1):5–19
- Ahipaşaoglu SD, Todd MJ (2013) A modified Frank–Wolfe algorithm for computing minimum-area enclosing ellipsoidal cylinders: Theory and algorithms. *Comput Geom* 46(5):494–519
- Allen-Zhu Z, Hazan E, Hu W, Li Y (2017) Linear convergence of a Frank–Wolfe type algorithm over trace-norm balls. *Adv Neural Inf Process Syst* 2017:6192–6201
- Bach F et al (2013) Learning with submodular functions: A convex optimization perspective. *Foundations and Trends®. Mach Learn* 6(2–3):145–373
- Bashiri MA, Zhang X (2017) Decomposition-invariant conditional gradient for general polytopes with line search. In: *Advances in neural information processing systems*, pp 2690–2700
- Beck A, Pauwels E, Sabach S (2015) The cyclic block conditional gradient method for convex optimization problems. *SIAM J Optim* 25(4):2024–2049
- Beck A, Shtern S (2017) Linearly convergent away-step conditional gradient for non-strongly convex functions. *Math Program* 164(1–2):1–27

- Berrada L, Zisserman A, Kumar MP (2018) Deep Frank–Wolfe for neural network optimization. In: International conference on learning representations
- Bertsekas DP (2015) Convex optimization algorithms. Athena Scientific, Nashua
- Bomze IM (1997) Evolution towards the maximum clique. *J Global Optim* 10(2):143–164
- Bomze IM, Budinich M, Pardalos PM, Pelillo M (1999) The maximum clique problem. In: Du D-Z, Pardalos P (eds) Handbook of combinatorial optimization, pp. 1–74. Springer
- Bomze IM, de Klerk E (2002) Solving standard quadratic optimization problems via linear, semidefinite and copositive programming. *J Global Optim* 24(2):163–185
- Bomze IM, Rinaldi F, Rota Bulò S (2019) First-order methods for the impatient: Support identification in finite time with convergent Frank–Wolfe variants. *SIAM J Optim* 29(3):2211–2226
- Bomze IM, Rinaldi F, Zeffiro D (2020) Active set complexity of the away-step Frank–Wolfe algorithm. *SIAM J Optim* 30(3):2470–2500
- Boyd S, Boyd SP, Vandenberghe L (2004) Convex optimization. Cambridge University Press, Cambridge
- Braun G, Pokutta S, Tu D, Wright S (2019) Blended conditional gradients. In: International conference on machine learning, PMLR, pp 735–743
- Braun G, Pokutta S, Zink D (2017) Lazifying conditional gradient algorithms. In: ICML, pp 566–575
- Candès EJ, Recht B (2009) Exact matrix completion via convex optimization. *Found Comput Math* 9(6):717–772
- Canon MD, Cullum CD (1968) A tight upper bound on the rate of convergence of Frank–Wolfe algorithm. *SIAM J Control* 6(4):509–516
- Carderera A, Pokutta S (2020) Second-order conditional gradient sliding. arXiv preprint [arXiv:2002.08907](https://arxiv.org/abs/2002.08907)
- Carlini N, Wagner D (2017) Towards evaluating the robustness of neural networks. In: 2017 IEEE symposium on security and privacy (sp), IEEE, pp 39–57
- Chakrabarty D, Jain P, Kothari P (2014) Provable submodular minimization using Wolfe’s algorithm. *Adv Neural Inform Process Syst* 27:802–809
- Chen J, Zhou D, Yi J, Gu Q (2020) A Frank–Wolfe framework for efficient and effective adversarial attacks. In: Proceedings of the AAAI conference on artificial intelligence vol 34, pp 3486–3494
- Chen PY, Zhang H, Sharma Y, Yi J, Hsieh CJ (2017) ZOO: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In: Proceedings of the 10th ACM workshop on artificial intelligence and security, pp 15–26
- Chen SS, Donoho DL, Saunders MA (2001) Atomic decomposition by basis pursuit. *SIAM Rev* 43(1):129–159
- Clarkson KL (2010) Coresets, sparse greedy approximation, and the Frank–Wolfe algorithm. *ACM Trans Algorithms* 6(4):1–30
- Combettes C, Pokutta S (2020) Boosting Frank–Wolfe by chasing gradients. In: International Conference on Machine Learning, PMLR, pp 2111–2121
- Combettes CW, Pokutta S (2021) Complexity of linear minimization and projection on some sets. arXiv preprint [arXiv:2101.10040](https://arxiv.org/abs/2101.10040)
- Cristofari A, De Santis M, Lucidi S, Rinaldi F (2020) An active-set algorithmic framework for non-convex optimization problems over the simplex. *Comput Optim Appl* 77:57–89
- Demyanov VF, Rubinov AM (1970) Approximate methods in optimization problems. American Elsevier, New York
- Devolder O, Glineur F, Nesterov Y (2014) First-order methods of smooth convex optimization with inexact oracle. *Math Program* 146(1):37–75
- Ding L, Fei Y, Xu Q, Yang C (2020) Spectral Frank–Wolfe algorithm: Strict complementarity and linear convergence. In: International conference on machine learning, PMLR, pp 2535–2544
- Dunn JC (1979) Rates of convergence for conditional gradient algorithms near singular and nonsingular extremals. *SIAM J Control Optim* 17(2):187–211
- Dunn JC, Harshbarger S (1978) Conditional gradient algorithms with open loop step size rules. *J Math Anal Appl* 62(2):432–444
- Ferreira O, Sosa W (2021) On the Frank–Wolfe algorithm for non-compact constrained optimization problems. *Optimization* 1–15
- Frank M, Wolfe P (1956) An algorithm for quadratic programming. *Naval Res Logist Q* 3(1–2):95–110
- Freund RM, Grigas P (2016) New analysis and results for the Frank–Wolfe method. *Math Program* 155(1–2):199–230
- Freund RM, Grigas P, Mazumder R (2017) An extended Frank–Wolfe method with in-face directions, and its application to low-rank matrix completion. *SIAM J Optim* 27(1):319–346

- Fujishige S (1980) Lexicographically optimal base of a polymatroid with respect to a weight vector. *Math Oper Res* 5(2):186–196
- Fukushima M (1984) A modified Frank–Wolfe algorithm for solving the traffic assignment problem. *Trans Res Part B Methodol* 18(2):169–177
- Garber D (2019) Linear convergence of Frank–Wolfe for rank-one matrix recovery without strong convexity. arXiv preprint [arXiv:1912.01467](https://arxiv.org/abs/1912.01467)
- Garber D (2020) Revisiting Frank–Wolfe for polytopes: Strict complementarity and sparsity. *Adv Neural Inform Process Syst* 33:18883–18893
- Garber D, Hazan E (2015) Faster rates for the Frank–Wolfe method over strongly-convex sets. *ICML* 15:541–549
- Garber D, Hazan E (2016) A linearly convergent variant of the conditional gradient algorithm under strong convexity, with applications to online and stochastic optimization. *SIAM J Optim* 26(3):1493–1528
- Goodfellow I, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, Courville A, Bengio Y (2014) Generative adversarial nets. In: *Advances in neural information processing systems*, pp 2672–2680
- Guélat J, Marcotte P (1986) Some comments on Wolfe’s away step. *Math Program* 35(1):110–119
- Gutman DH, Pena JF (2021) The condition number of a function relative to a set. *Math Program* 188:255–294
- Harchaoui Z, Juditsky A, Nemirovski A (2015) Conditional gradient algorithms for norm-regularized smooth convex optimization. *Math Program* 152(1):75–112
- Hogan WW (1971) Convergence results for some extensions of the Frank–Wolfe method. Tech. rep., California Univ Los Angeles Western Management Science Inst
- Holloway CA (1974) An extension of the Frank and Wolfe method of feasible directions. *Math Program* 6(1):14–27
- Hungerford JT, Rinaldi F (2019) A general regularized continuous formulation for the maximum clique problem. *Math Oper Res* 44(4):1161–1173
- Jaggi M (2011) Sparse convex optimization methods for machine learning. Ph.D. thesis, ETH Zurich
- Jaggi M (2013) Revisiting Frank–Wolfe: Projection-free sparse convex optimization. *ICML* 1:427–435
- Jaggi M, Sulovský M (2010) A simple algorithm for nuclear norm regularized problems. In: *ICML*, pp 471–478
- Joulin A, Tang K, Fei-Fei L (2014) Efficient image and video co-localization with Frank–Wolfe algorithm. In: *European conference on computer vision*. Springer, pp 253–268
- Kazemi E, Kerdreux T, Wang L (2021) Generating structured adversarial attacks using Frank–Wolfe method. arXiv preprint [arXiv:2102.07360](https://arxiv.org/abs/2102.07360)
- Kerdreux T, d’Aspremont A, Pokutta S (2021) Projection-free optimization on uniformly convex sets. In: *International Conference on Artificial Intelligence and Statistics*, pp. 19–27. PMLR
- Kerdreux T, Liu L, Lacoste-Julien S, Scieur D (2020) Affine invariant analysis of Frank–Wolfe on strongly convex sets. arXiv preprint [arXiv:2011.03351](https://arxiv.org/abs/2011.03351)
- Konnov I (2018) Simplified versions of the conditional gradient method. *Optimization* 67(12):2275–2290
- Kumar P, Mitchell JS, Yıldırım EA (2003) Approximate minimum enclosing balls in high dimensions using core-sets. *J Exp Algorithmics* 8:1–1
- Lacoste-Julien S (2016) Convergence rate of Frank–Wolfe for non-convex objectives. arXiv preprint [arXiv:1607.00345](https://arxiv.org/abs/1607.00345)
- Lacoste-Julien S, Jaggi M (2015) On the global linear convergence of Frank–Wolfe optimization variants. In: *Advances in neural information processing systems*, pp 496–504
- Lacoste-Julien S, Jaggi M, Schmidt M, Pletscher P (2013) Block-coordinate Frank–Wolfe optimization for structural SVMs. In: Dasgupta S, McAllester D (eds) *Proceedings of the 30th International Conference on Machine Learning, Proceedings of Machine Learning Research*, vol 28, PMLR, Atlanta, Georgia, USA, pp 53–61
- Lan G (2020) *First-order and stochastic optimization methods for machine learning*. Springer, New York
- Lan G, Zhou Y (2016) Conditional gradient sliding for convex optimization. *SIAM J Optim* 26(2):1379–1409
- LeBlanc LJ, Morlok EK, Pierskalla WP (1975) An efficient approach to solving the road network equilibrium traffic assignment problem. *Transp Res* 9(5):309–318
- Levitin ES, Polyak BT (1966) Constrained minimization methods. *USSR Comput Math Math Phys* 6(5):1–50
- Locatello F, Khanna R, Tschannen M, Jaggi M (2017) A unified optimization view on generalized matching pursuit and Frank–Wolfe. In: *Artificial intelligence and statistics*. PMLR, pp 860–868

- Luce RD, Perry AD (1949) A method of matrix analysis of group structure. *Psychometrika* 14(2):95–116
- Mangasarian O (1996) Machine learning via polyhedral concave minimization. *Appl Math Parallel Comput*. Springer, New York, pp 175–188
- Mitchell B, Demyanov VF, Malozemov V (1974) Finding the point of a polyhedron closest to the origin. *SIAM J Control* 12(1):19–26
- Mitradjieva M, Lindberg PO (2013) The stiff is moving–conjugate direction Frank–Wolfe methods with applications to traffic assignment. *Transp Sci* 47(2):280–293
- Mu C, Zhang Y, Wright J, Goldfarb D (2016) Scalable robust matrix recovery: Frank–Wolfe meets proximal methods. *SIAM J Sci Comput* 38(5):A3291–A3317
- Osokin A, Alayrac JB, Lukasewitz I, Dokania P, Lacoste-Julien S (2016) Minding the gaps for block Frank–Wolfe optimization of structured svms. In: *International conference on machine learning*, PMLR, pp 593–602
- Peña J, Rodriguez D (2018) Polytope conditioning and linear convergence of the Frank–Wolfe algorithm. *Math Oper Res* 44(1):1–18
- Pedregosa F, Negiar G, Askari A, Jaggi M (2020) Linearly convergent Frank–Wolfe with backtracking line-search. In: *International conference on artificial intelligence and statistics*. PMLR, pp 1–10
- Perederieieva O, Ehr Gott M, Raith A, Wang JY (2015) A framework for and empirical study of algorithms for traffic assignment. *Comput Oper Res* 54:90–107
- Rademacher L, Shu C (2020) The smoothed complexity of Frank–Wolfe methods via conditioning of random matrices and polytopes. *arXiv preprint arXiv:2009.12685*
- Rinaldi F, Schoen F, Sciandrone M (2010) Concave programming for minimizing the zero-norm over polyhedral sets. *Comput Optim Appl* 46(3):467–486
- Rinaldi F, Zeffiro D (2020) Avoiding bad steps in Frank Wolfe variants. *arXiv preprint arXiv:2012.12737*
- Rinaldi F, Zeffiro D (2020) A unifying framework for the analysis of projection-free first-order methods under a sufficient slope condition. *arXiv preprint arXiv:2008.09781*
- Sahu AK, Kar S (2020) Decentralized zeroth-order constrained stochastic optimization algorithms: Frank–Wolfe and variants with applications to black-box adversarial attacks. *Proc IEEE* 108(11):1890–1905
- Shah N, Kolmogorov V, Lampert CH (2015) A multi-plane block-coordinate Frank–Wolfe algorithm for training structural svms with a costly max-oracle. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp 2737–2745
- Sun Y (2020) Safe screening for the generalized conditional gradient method. *Image* 1:2
- Tibshirani R (1996) Regression shrinkage and selection via the lasso. *J Roy Stat Soc: Ser B (Methodol)* 58(1):267–288
- Vapnik V (2013) *The nature of statistical learning theory*. Springer, New York
- Von Hohenbalken B (1977) Simplicial decomposition in nonlinear programming algorithms. *Math Program* 13(1):49–68
- Wang H, Lu H, Mazumder R (2020) Frank–Wolfe methods with an unbounded feasible region and applications to structured learning. *arXiv preprint arXiv:2012.15361*
- Wang YX, Sadhanala V, Dai W, Neiswanger W, Sra S, Xing E (2016) Parallel and distributed block-coordinate Frank–Wolfe algorithms. In: *International Conference on Machine Learning*. PMLR, pp 1548–1557
- Wardrop JG (1952) Road paper. some theoretical aspects of road traffic research. *Proc Inst Civ Eng* 1(3):325–362
- Weintraub A, Ortiz C, González J (1985) Accelerating convergence of the Frank–Wolfe algorithm. *Transp Res Part B Methodol* 19(2):113–122
- Wolfe P (1970) Convergence theory in nonlinear programming. In: Abadie J (ed) *Integer and nonlinear programming*. North Holland, pp 1–36
- Wolfe P (1976) Finding the nearest point in a polytope. *Math Program* 11(1):128–149
- Wu Q, Hao JK (2015) A review on algorithms for maximum clique problems. *Eur J Oper Res* 242(3):693–709
- Xu Y, Yang T (2018) Frank–Wolfe method is automatically adaptive to error bound condition. *arXiv preprint arXiv:1810.04765*
- Yıldırım EA (2008) Two algorithms for the minimum enclosing ball problem. *SIAM J Optim* 19(3):1368–1391