
Step-Size Adaptivity in Projection-Free Optimization

Fabian Pedregosa
ETH Zürich / UC Berkeley

Armin Askari, Geoffrey Negiar
UC Berkeley

Martin Jaggi
EPFL

Abstract

We propose and analyze adaptive step-size variants of the Frank-Wolfe, its Away-Steps and Pairwise variants as well as Matching Pursuit. The proposed methods leverage local information of the objective through a backtracking line search strategy. This has two key advantages: First, it does not require to estimate problem-dependent constants that might be costly to compute, such as the Lipschitz or the curvature constant. Second, the proposed criterion is adaptive to local properties of the objective, allowing for larger step-sizes. For all proposed methods, we derive convergence rates on convex and non-convex objectives that asymptotically match the strongest known bounds for non-adaptive variants. As a side-product of this analysis, we obtain the first linear convergence rate of Frank-Wolfe without exact line search and the first known bounds for matching pursuit on non-convex objectives. Benchmarks on three different datasets illustrate the computational advantage of the proposed methods.

1 Introduction

Most first-order optimization methods rely on a step-size parameter. It controls the magnitude of the update and has a crucial impact on its performance: a step-size that is too small will give an unnecessary slow convergence, while a step-size that is too big might lead to divergence. Different techniques have been developed to estimate this parameter, among which:

Exact line search. In very specific cases it is possible to derive an analytical solution for the step-size that minimizes the objective along the update direction.

This technique gives the best performance but is typically only available for quadratic loss functions.

Exact line search on quadratic upper bound. This approach is more widely applicable but requires knowledge of global constants of the objective like the gradient's Lipschitz constant. These can be difficult to estimate and are often conservative, leading to suboptimal step sizes.

Adaptive step size or backtracking line search instead selects a step-size based on an implicit equation that guarantees sufficient progress is made at each iteration. Unlike the previous approach, it does not require knowledge of global constants and typically result in larger step sizes. Adaptive step size methods are classical in the context of smooth unconstrained optimization (Goldstein, 1965; Armijo, 1966; Wolfe, 1969) and have also been extended to projection (or proximal)-based methods (Beck and Teboulle, 2009; Nesterov, 2013).

The Frank-Wolfe (FW) or conditional gradient algorithm (Frank and Wolfe, 1956; Demyanov and Rubinov, 1967) is one of the oldest methods for non-linear constrained optimization and has experienced a renewed interest in recent years due to its applications in machine learning (Jaggi, 2013). Rather surprisingly, and with a few exceptions that we revisit in §4, adaptive step size variants remain largely unexplored. Furthermore, variants like Away-Steps FW and Pairwise FW have recently been shown to enjoy linear convergence over a polytope domain (Lacoste-Julien and Jaggi, 2015) but rely on exact line search, severely limiting their applicability.

This raises our motivating question: *Is it possible to develop adaptive step size variants of FW that achieve the same rate of convergence as those that rely on exact line search?* In this paper we give a positive answer and extend the method to other projection-free methods such as Matching Pursuit (MP) (Mallat and Zhang, 1993).

Outline and main contributions. Our main contribution is the design and analysis of adaptive step-size variants of FW, Away-steps FW, Pairwise FW and

MP. In all cases, we develop a convergence rate analysis that matches the best known bounds on convex, strongly convex and non-convex problems, including linear convergence over polytopes for Away-Step FW and Pairwise FW. In the case of MP, we provide the first convergence rates for non-convex objectives. The paper is structured as follows:

- *Methods.* In §2 we describe the Adaptive FW the linearly-convergent variants Away-steps FW and Pairwise FW. §3 describes an adaptive variant of MP. In §4 we relate the proposed methods with previous literature.
- *Analysis.* §5 provides a convergence analysis of the proposed methods on non-convex, convex and strongly convex objectives.
- *Experiments.* §6 compares the proposed methods against existing approaches on 3 different datasets and 3 problems.

Notation. Throughout the paper we denote vectors and vector-valued functions in lowercase boldface (e.g. \mathbf{x} or $\mathbf{arg\,min}$), matrices in uppercase boldface letters (e.g. \mathbf{D}), and sets in caligraphic letters (e.g., \mathcal{A}). We say a function f is L -smooth if it is differentiable and its gradient is L -Lipschitz continuous, that is, if it verifies $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$ for all \mathbf{x}, \mathbf{y} in the domain. A function is μ -strongly convex if $f - \frac{\mu}{2}\|\cdot\|^2$ is convex.

2 Adaptive Frank-Wolfe

In this section we consider constrained optimization problems of the form

$$\underset{\mathbf{x} \in \text{conv}(\mathcal{A})}{\text{minimize}} \quad f(\mathbf{x}), \quad (\text{OPT-FW})$$

where f is L -smooth and $\text{conv}(\mathcal{A})$ is the convex hull of a potentially infinite but bounded set of elements (which we will refer to as *atoms*) in \mathbb{R}^p .

The proposed adaptive FW variant (named **AdaFW**) is specified in Algorithm 1. As FW, it requires to solve a linear subproblem of the form $\mathbf{arg\,min}_{\mathbf{s} \in \mathcal{A}} \langle \nabla f(\mathbf{x}_t), \mathbf{s} \rangle$ at each iteration. In practice, however, these are often not solved exactly but only to a predetermined accuracy. We explicitly allow it through an optional *subproblem quality* parameter $\delta \in (0, 1]$ and consider in line 3 the linear subproblem of finding $\mathbf{s}_t \in \mathcal{A}$ such that

$$\langle \nabla f(\mathbf{x}_t), \mathbf{s}_t - \mathbf{x}_t \rangle \leq \delta \min_{\mathbf{s} \in \mathcal{A}} \langle \nabla f(\mathbf{x}_t), \mathbf{s} - \mathbf{x}_t \rangle. \quad (1)$$

When $\delta = 1$ the above is equivalent to solving the linear subproblems exactly, but for $\delta < 1$ the solutions are

Algorithm 1: Adaptive Frank-Wolfe (AdaFW)

```

1 Input:  $\mathbf{x}_0 \in \text{conv}(\mathcal{A})$ , initial Lipschitz estimate
    $L_{-1} > 0$ , algorithm tolerance  $\varepsilon \geq 0$ , subproblem
   quality  $\delta \in (0, 1]$ , adaptivity params  $\tau > 1, \eta \geq 1$ 
2 for  $t = 0, 1 \dots$  do
3   Choose any  $\mathbf{s}_t \in \mathcal{A}$  that verifies (1)
4   Set  $\mathbf{d}_t = \mathbf{s}_t - \mathbf{x}_t$  and  $g_t = \langle -\nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle$ 
5   if  $g_t \leq \delta \varepsilon$  then return  $\mathbf{x}_t$ ;
6   Set  $M = L_{t-1}/\eta$ ,  $\gamma = \min \{g_t/(M\|\mathbf{d}_t\|^2), 1\}$ 
7   while  $f(\mathbf{x}_t + \gamma \mathbf{d}_t) > Q_t(\gamma, M)$  do
8      $M = \tau M$ ,  $\gamma = \min \{g_t/(M\|\mathbf{d}_t\|^2), 1\}$ 
9   Set  $L_t = M$  and  $\gamma_t = \gamma$ 
10  Set  $\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma_t \mathbf{d}_t$ 

```

allowed to be suboptimal. This is the same criterion as in Lacoste-Julien et al. (2013); Locatello et al. (2017).

The (approximate) solution to the linear subproblems gives the update direction $\mathbf{d}_t = \mathbf{s}_t - \mathbf{x}_t$, so that the next iterate is of the form $\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma_t \mathbf{d}_t$. To set the step-size γ_t we consider the function

$$Q_t(\gamma, M) = f(\mathbf{x}_t) + \gamma \langle \nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle + \frac{\gamma^2 M}{2} \|\mathbf{d}_t\|^2. \quad (2)$$

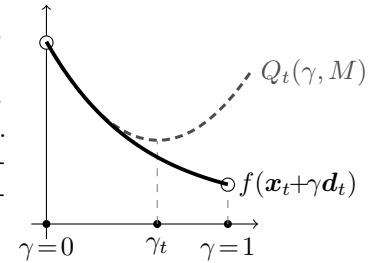
This function is quadratic in its first argument, and minimizing for this parameter in the interval $[0, 1]$ gives $\gamma_t = \min \{g_t/(M\|\mathbf{d}_t\|^2), 1\}$, which is the step size used in lines 6 and 8.

This surrogate depends on a scalar parameter $M \geq 0$ that we will refer to as the *Lipschitz estimate*. This parameter is initialized as $M = L_{t-1}/\eta$, where $\eta \geq 1$ allows it to decrease along iterations. Lines 7–8 then increase this estimate by a power factor of $\tau > 1$ until the following condition is verified:

$$f(\mathbf{x}_t + \gamma \mathbf{d}_t) \leq Q_t(\gamma, M), \quad \gamma = \min \{g_t/(M\|\mathbf{d}_t\|^2), 1\}.$$

We call this condition the *sufficient decrease condition*.

Once this condition is verified, the current step-size is accepted and the value of M is assigned the name L_t . Geometrically, the sufficient decrease condition ensures that the quadratic surrogate



$Q_t(\cdot, M)$ at its constrained minimum γ_t is an upper bound of $\gamma \mapsto f(\mathbf{x}_t + \gamma \mathbf{d}_t)$. We emphasize that unlike the “exact line search on quadratic upper bound” approach, in this case the surrogate Q_t need not be a global upper bound on the objective.

Overhead of the adaptive step-size strategy. Evaluation of the sufficient decrease condition requires

two extra evaluations of the objective function. If the condition verified, then it is only evaluated at the current and next iterate. FW requires anyway to compute the gradient at these iterates, hence in cases in which the objective function is available as a byproduct of the the gradient this overhead becomes negligible.

Furthermore, it is possible to bound the total number of evaluation of the sufficient decrease condition (Lemma 4, Appendix B) by

$$\left[1 + \frac{\log \eta}{\log \tau}\right] (t+1) + \frac{1}{\log \tau} \max \left\{ \log \frac{\tau L}{L_{-1}}, 0 \right\}. \quad (3)$$

For the adaptive step size parameters we recommend $\eta = 1.001$, $\tau = 2$. Using these values and for $L_{-1} \geq L/10$ the above bounds the number of evaluation by $\approx 1.001(t+1) + 4.32$. This implies that for $t \geq 1000$, 99% of the iterations will only perform one evaluation of the sufficient decrease condition.

The proposed algorithm depends also on an initial value for the Lipschitz estimate L_{-1} . A simple heuristic which we use in the experiments consists in starting with $L_{-1} = 10^{-3}$, $\tilde{\mathbf{x}} = \mathbf{x}_0 - (1/L_{-1})\nabla f(\mathbf{x}_0)$ and multiply L_{-1} by 10 until $f(\tilde{\mathbf{x}}) \leq f(\mathbf{x}_0)$.

2.1 Away-Steps and Pairwise Variants

In this subsection we present adaptive step-size variants of the Away-Steps FW (named **AdaAFW**) and Pairwise FW (named **AdaPFW**). The Away-Steps FW (Guélat and Marcotte, 1986) is a popular variant of FW that adds the option to move away from an atom in the current representation of the iterate. In the case of a polytope domain, it was recently shown to enjoy a linear convergence rate for strongly convex objectives (Garber and Hazan, 2013; Lacoste-Julien and Jaggi, 2015). The Pairwise FW was proposed by Lacoste-Julien and Jaggi (2015) based on the MDM algorithm of Mitchell et al. (1974).

Unlike **Adaptive FW**, the methods introduced in this subsection require to keep track of previous updates. For this purpose we introduce the *active set* $\mathcal{S}_t \subseteq \mathcal{A}$, which contains the atoms \mathbf{s}_t selected by previous iterations that have non-zero weight $\alpha_{\mathbf{s},t} > 0$ in the expansion $\mathbf{x}_t = \sum_{\mathbf{s} \in \mathcal{S}_t} \alpha_{\mathbf{s},t} \mathbf{s}$.

Just like their non-adaptive variants, both algorithms must solve two linear subproblems per iteration. The first one (Line 4) is identical to that **AdaFW**. The second one is different, and is restricted to elements in the active set. More precisely, it consists in finding any $\mathbf{s}_t \in \mathcal{S}_t$ that verifies

$$\langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{v}_t \rangle \leq \delta \min_{\mathbf{v} \in \mathcal{S}_t} \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{v} \rangle. \quad (4)$$

Note that this subproblem can also be solved approximately through a quality parameter $\delta \in (0, 1]$.

Alg. 2: Adaptive Away-Steps FW (AdaAFW)

```

1  $\mathbf{x}_0 \in \mathcal{A}$ , initial Lipschitz estimate  $L_{-1} > 0$ , tolerance
    $\varepsilon \geq 0$ , subproblem quality  $\delta \in (0, 1]$ , adaptivity
   params  $\tau > 1, \eta \geq 1$ 
2 Let  $\mathcal{S}_0 = \{\mathbf{x}_0\}$  and  $\alpha_{0,\mathbf{v}} = 1$  for  $\mathbf{v} = \mathbf{x}_0$  and  $\alpha_{0,\mathbf{v}} = 0$ 
   otherwise.
3 for  $t = 0, 1 \dots$  do
4     Choose any  $\mathbf{s}_t \in \mathcal{A}$  that satisfies (1)
5     Choose any  $\mathbf{v}_t \in \mathcal{S}_t$  that satisfies (4)
6     if  $\langle \nabla f(\mathbf{x}_t), \mathbf{s}_t - \mathbf{x}_t \rangle \leq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{v}_t \rangle$  then
7          $\mathbf{d}_t = \mathbf{s}_t - \mathbf{x}_t$  and  $\gamma_t^{\max} = 1$ 
8     else
9          $\mathbf{d}_t = \mathbf{x}_t - \mathbf{v}_t$ , and  $\gamma_t^{\max} = \alpha_{\mathbf{v}_t,t} / (1 - \alpha_{\mathbf{v}_t,t})$ 
10    Set  $g_t = \langle -\nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle$ 
11    if  $g_t \leq \delta \varepsilon$  then return  $\mathbf{x}_t$ ;
12    Set  $M = L_{t-1}/\eta$ ,  $\gamma = \min \{g_t / (M \|\mathbf{d}_t\|^2), \gamma_{\max}\}$ 
13    while  $f(\mathbf{x}_t + \gamma \mathbf{d}_t) > Q_t(\gamma, M)$  do
14         $M = \tau M$ ,  $\gamma = \min \{g_t / (M \|\mathbf{d}_t\|^2), \gamma_{\max}\}$ 
15    Set  $L_t = M$  and  $\gamma_t = \gamma$ 
16     $\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma_t \mathbf{d}_t$ 
17    Update active set  $\mathcal{S}_{t+1}$  and  $\alpha_{t+1}$  (see text)
```

Algorithm 3: Adaptive Pairwise FW (AdaPFW)

```

1 As AdaAFW, replacing Lines 6–9 by
2  $\mathbf{d}_t = \mathbf{s}_t - \mathbf{v}_t$  and  $\gamma_t^{\max} = \alpha_{\mathbf{v}_t,t}$ 
```

For **AdaAFW** we construct two potential descent directions: the FW direction $\mathbf{s}_t - \mathbf{x}_t$ and the Away direction $\mathbf{x}_t - \mathbf{v}_t$. The chosen direction is the one that correlates the most with the negative gradient. **AdaPFW** instead constructs a single descent direction $\mathbf{d}_t = \mathbf{s}_t - \mathbf{v}_t$ using the result of both linear subproblems.

Lines 12–14 sets the Lipschitz estimate. This is done using a sufficient decrease similar to the one of **AdaFW** but with the maximum step-size of 1 replaced by γ_t^{\max} , the latter being set to ensure that the iterates remain within the domain.

Updating the support and associated coefficients. For **AdaAFW** these can be updated as follows. In case of a FW step we update the support set $\mathcal{S}_{t+1} = \{\mathbf{s}_t\}$ if $\gamma_t = 1$ and otherwise $\mathcal{S}_{t+1} = \mathcal{S}_t \cup \{\mathbf{s}_t\}$, with coefficients $\alpha_{\mathbf{v},t+1} = (1 - \gamma_t)\alpha_{\mathbf{v},t}$ for $\mathbf{v} \in \mathcal{S}_t \setminus \{\mathbf{s}_t\}$ and $\alpha_{\mathbf{s}_t,t+1} = (1 - \gamma_t)\alpha_{\mathbf{s}_t,t} + \gamma_t$. In case of an away step we instead have the following update rule: If $\gamma_t = \gamma_{\max}$, then $\mathcal{S}_{t+1} = \mathcal{S}_t \setminus \{\mathbf{v}_t\}$, and if $\gamma_t < \gamma_{\max}$, then $\mathcal{S}_{t+1} = \mathcal{S}_t$. Finally, we update the weights as $\alpha_{\mathbf{v},t+1} = (1 + \gamma_t)\alpha_{\mathbf{v},t}$ for $\mathbf{v} \in \mathcal{S}_t \setminus \{\mathbf{v}_t\}$ and $\alpha_{\mathbf{v}_t,t+1} = (1 + \gamma_t)\alpha_{\mathbf{v}_t,t} - \gamma_t$ for the other atoms.

AdaPFW on the other hand only moves weight from \mathbf{v}_t to \mathbf{s}_t and so the update for the coefficients becomes

$\alpha_{s_t, t+1} = \alpha_{s_t, t} + \gamma_t$, $\alpha_{v_t, t+1} = \alpha_{v_t, t} - \gamma_t$, with $\mathcal{S}_{t+1} = (\mathcal{S}_t \setminus \{v_t\}) \cup \{s_t\}$ if $\alpha_{v_t, t+1} = 0$ and $\mathcal{S}_{t+1} = \mathcal{S}_t \cup \{s_t\}$ otherwise.

3 Adaptive Matching Pursuit

Matching Pursuit (Mallat and Zhang, 1993; Locatello et al., 2017) is an algorithm to solve optimization problems of the form

$$\underset{\mathbf{x} \in \text{lin}(\mathcal{A})}{\text{minimize}} \ f(\mathbf{x}), \quad (\text{OPT-MP})$$

where $\text{lin}(\mathcal{A}) \stackrel{\text{def}}{=} \{\sum_{v \in \mathcal{A}} \lambda_v v \mid \lambda_v \in \mathbb{R}\}$ is the linear span of the set of atoms \mathcal{A} . As for the Adaptive FW algorithm, we assume that f is L -smooth and \mathcal{A} a potentially infinite but bounded set of elements in \mathbb{R}^p .

The MP algorithm relies on solving at each iteration a linear subproblem over the set $\mathcal{B} \stackrel{\text{def}}{=} \mathcal{A} \cup -\mathcal{A}$, with $-\mathcal{A} = \{-\mathbf{a} \mid \mathbf{a} \in \mathcal{A}\}$. The linear subproblem that needs to be solved at each iteration is the following, where as for previous variants, we allow for an optional quality parameter $\delta \in (0, 1]$:

$$\langle \nabla f(\mathbf{x}_t), \mathbf{s}_t \rangle \leq \delta \min_{\mathbf{s} \in \mathcal{B}} \langle \nabla f(\mathbf{x}_t), \mathbf{s} \rangle. \quad (5)$$

In Algorithm 4 we detail a novel adaptive variant of the MP algorithm, which we name **AdaMP**. It relies on a Lipschitz estimate M verifying sufficient decrease condition (Lines 7–8)

$$f(\mathbf{x}_t + \gamma \mathbf{d}_t) \leq Q_t(\gamma, M) \text{ with } \gamma = g_t / (M \|\mathbf{d}_t\|^2) \quad (6)$$

Note that unlike previous approaches, the step-size is unconstrained in this case.

4 Related work

We comment on the most closely related ideas, summarized in Table 1. Adaptive step size variants of FW have been described in (Dunn, 1980) and (Beck et al., 2015).

Dunn (1980) adapted the Goldstein (1965) and Armijo (1966) backtracking line search methods to select the step size in the FW method. For the Armijo criterion this method requires to select parameters $\eta \in (0, 1)$ and $\delta \in (0, \frac{1}{2}]$. Then the step size is chosen as $\gamma_t = \eta^i$, where i is the smallest integer such that

$$f(\mathbf{x}_t + \gamma_t \mathbf{d}_t) \leq f(\mathbf{x}_t) - \delta \gamma_t g_t. \quad (7)$$

A crucial difference with our approach is that here there is no parameter that estimates the smoothness of the objective like the Lipschitz estimate L_t . Since

Alg. 4: Adaptive Matching Pursuit (AdaMP)

```

1  $\mathbf{x}_0 \in \mathcal{D}$ , initial Lipschitz estimate  $L_{-1} > 0$ , tolerance
    $\varepsilon \geq 0$ , subproblems quality  $\delta \in (0, 1]$ , adaptivity
   params  $\tau > 1, \eta \geq 1$ 
2 for  $t = 0, 1 \dots$  do
3   Choose any  $\mathbf{s}_t \in \mathcal{B}$  that verifies (5)
4   Set  $\mathbf{d}_t = \mathbf{s}_t$  and  $g_t = \langle -\nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle$ 
5   if  $g_t \leq \delta \varepsilon$  then return  $\mathbf{x}_t$ ;
6   Set  $M = L_{t-1}/\eta$ ,  $\gamma = g_t / (M \|\mathbf{d}_t\|^2)$ 
7   while  $f(\mathbf{x}_t + \gamma \mathbf{d}_t) > Q_t(\gamma, M)$  do
8      $M = \tau M$ ,  $\gamma = g_t / (M \|\mathbf{d}_t\|^2)$ 
9   Set  $L_t = M$  and  $\gamma_t = \gamma$ 
10   $\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma_t \mathbf{d}_t$ 
```

the update in FW involves an extremal element that changes from one iterate to the next, the step size can also change drastically. In practice this leads to a much larger number of line search corrections. Furthermore, there is no bound on the number of evaluations of its sufficient decrease condition and it has not been extended to the linearly-convergent FW variants.

Beck et al. (2015) proposed a different adaptive FW variant for a cyclical variant of the block-coordinate FW (Lacoste-Julien et al., 2013). For the case of a single block of variables, it gives an adaptive variant of FW. In this case, the step size is set as $\gamma = \min\{g_t / (M_t \|\mathbf{d}_t\|^2), 1\}$, where M_t is selected such that the following is verified:

$$f(\mathbf{x}_t + \gamma_t \mathbf{d}_t) \leq f(\mathbf{x}_t) - \frac{\gamma}{2} g_t. \quad (8)$$

Since $g_t \geq \gamma_t L_t \|\mathbf{x}_t - \mathbf{s}_t\|^2$ (Beck et al., 2015, Lemma 4.6), our sufficient decrease condition implies the above and so leads to larger step sizes. Another aspect of practical importance is that in this algorithm the Lipschitz estimates M_t are required to be monotonically increasing, in contrast with the proposed methods in which the Lipschitz estimates are also allowed to decrease. Furthermore, in (Beck et al., 2015) there is no explicit bound on the number of evaluations of the sufficient decrease condition, no convergence guarantees for non-convex objectives, nor any extension to linearly-convergent FW variants.

In the context of MP, the vast majority of existing literature assumes a quadratic objective and consequently access to an exact line search. Locatello et al. (2017) recently proposed a variant of MP that only requires the objective to be smooth (instead of quadratic). In this algorithm the step-size depends on the Lipschitz constant of ∇f . We propose an adaptive variant of this algorithm in which the global Lipschitz constant is replaced by a local estimate.

The proposed **AdaMP** algorithm is most similar to the

	Related work	non-convex analysis	approximate subproblems	linear convergence	adaptive step size	bounded backtracking
Frank-Wolfe	<i>This work</i>	✓	✓	✓	✓	✓
	(Lacoste-Julien and Jaggi, 2015)	✗	✗	✓	✗	N/A
	(Beck et al., 2015)	✗	✓ [†]	✗	✓	✗
	(Dunn, 1980)	✓	✗	✗	✓	✗
MP	<i>This work</i>	✓	✓	✓	✓	✓
	(Locatello et al., 2017)	✗	✓	✓	✗	N/A

Table 1: **Comparison with related work.** *non-convex analysis*: convergence guarantees for problems with a non-convex objective. *approximate subproblems*: convergence guarantees cover the case in which linear subproblems are solved approximately. *linear convergence*: guaranteed linear rate of convergence (under hypothesis). *bounded backtracking*: explicit bound for the total number of inner iterations in adaptive step size methods. [†] = assumes cartesian product structure of the domain

“Norm-Corrective Generalized Matching Pursuit” in (Locatello et al., 2017, Algorithm 4, variant 0). The difference between both algorithms lies in the choice of step-size: the variant of Locatello et al. (2017) uses the step-size $g_t L^{-1} \|\mathbf{d}_t\|^{-2}$, which relies on knowledge of the global Lipschitz constant L , while the proposed variant replaces it by the potentially much smaller and adaptive estimate L_t . Furthermore, the analysis of Locatello et al. (2017) was not extended to non-convex objectives.

5 Analysis

In this section, we provide a convergence rate analysis of the proposed methods, showing that all proposed methods enjoy a $\mathcal{O}(1/\sqrt{t})$ convergence rate for non-convex objectives (Theorem 1), a stronger $\mathcal{O}(1/t)$ convergence rate for convex objectives (Theorem 2), and linear convergence for strongly convex objectives for some algorithms and domains (Theorem 4).

Notation. In this section we make use of the following extra notation:

- We will use \mathcal{D} to refer to the domain of the function, which will be $\text{conv}(\mathcal{A})$ when referring to FW and $\text{lin}(\mathcal{A})$ when referring to MP.
- We denote the *objective suboptimality* at the t -th iteration as $h_t = f(\mathbf{x}_t) - \min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x})$.
- *Good and bad steps.* Our analysis, as that of Lacoste-Julien and Jaggi (2015), relies on a notion of “good” and “bad” steps. We define bad steps as those that verify $\gamma_t = \gamma_t^{\max}$ and $\gamma_t^{\max} < 1$ and good steps as any step that is not a bad step. Their name comes from the fact that for bad steps we will not be able to provide a meaningful decrease bound. Some algorithms do not have bad steps, e.g. AdaFW and AdaMP, while other algorithms do but allow their number to be bounded.

In order to provide a unified analysis of the different FW variants, we introduce the following notation. We denote by N_t the number of “good steps” up to iteration t . Crucially, it is possible to lower bound the number of good step for all algorithms as follows:

$$N_t = t \text{ for AdaFW and AdaMP,} \quad (9)$$

$$N_t \geq t/2 \text{ for AdaAFW,} \quad (10)$$

$$N_t \geq t/(3|\mathcal{A}| + 1) \text{ for AdaPFW} \quad (11)$$

where it is worth noting that the last bound for AdaPFW requires the set of atoms \mathcal{A} to be finite. The proof of these bounds can be found in Appendix B.1 and are a direct translation of those in (Lacoste-Julien and Jaggi, 2015). We have found these bounds to be very loose, as in practice the fraction of bad/good steps is negligible, commonly of the order of 10^{-5} (see last column of the table in Figure 1).

- *Average and maximum of Lipschitz estimate.* In order to highlight the better convergence rates that can be obtained by adaptive methods we introduce the average and maximum estimate over good step-sizes. Let \mathcal{G}_t denote the indices of good steps up to iteration t . Then we define the average and maximum Lipschitz estimate as

$$\bar{L}_t \stackrel{\text{def}}{=} \frac{1}{N_t} \sum_{k \in \mathcal{G}_t} L_k \quad (12)$$

$$L_t^{\max} \stackrel{\text{def}}{=} \max_{k \in \mathcal{G}_t} L_k \quad (13)$$

respectively. In the worse case, both quantities can be upper bounded by $\max\{\tau L, L_{-1}\}$ (Proposition 2), which can be used to obtain asymptotic convergence rates. This bound is however very pessimistic. We have found that in practice \bar{L}_t is often more than 100 times smaller than L (see second to last column of the table in Figure 1).

Our new convergence rates are presented in the following theorems, which consider the cases of non-convex,

convex and strongly convex objectives. The results are discussed in §5.4 and the proofs can be found in Appendix C, Appendix D and Appendix E respectively.

5.1 Non-convex objectives

Gap function. Convergence rates for convex and strongly convex functions are given in terms of the objective function suboptimality or a primal-dual gap. As the gap upper-bounds (i.e. certifies) the suboptimality, the latter is a stronger result in this scenario. In the case of non-convex objectives, as is common for first order methods, we will only be able to guarantee convergence to a stationary point, defined as any element $\mathbf{x}^* \in \mathcal{D}$ such that $\langle \nabla f(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle \geq 0$ for all $\mathbf{x} \in \mathcal{D}$ (Bertsekas, 1999).

Following Lacoste-Julien (2016); Reddi et al. (2016), for FW variants we will use as convergence criterion the FW gap, defined as $g^{\text{FW}}(\mathbf{x}) = \max_{\mathbf{s} \in \mathcal{D}} \langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{s} \rangle$. From the definition of stationary point it is clear that the FW gap is zero only at a stationary point.

In the context of MP, we propose the following criterion which we name the MP gap: $g^{\text{MP}}(\mathbf{x}) = \max_{\mathbf{s} \in \mathcal{B}} \langle \nabla f(\mathbf{x}), \mathbf{s} \rangle$. Note that g^{MP} is always non-negative and $g^{\text{MP}}(\mathbf{x}^*) = 0$ implies $\langle \nabla f(\mathbf{x}^*), \mathbf{s} \rangle = 0$ for all $\mathbf{s} \in \mathcal{B}$. By linearity of the inner product we then have $\langle \nabla f(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle = 0$ for any \mathbf{x} in the domain, since $\mathbf{x} - \mathbf{x}^*$ lies in the linear span of \mathcal{A} . Hence \mathbf{x}^* is a stationary point and g^{MP} is an appropriate measure of stationarity for this problem.

Theorem 1. *Let \mathbf{x}_t denote the iterate generated by any of the proposed algorithms after t iterations, with $N_{t+1} \geq 1$. Then we have:*

$$\lim_{t \rightarrow \infty} g(\mathbf{x}_t) = 0 \quad \text{and} \quad (14)$$

$$\min_{k=0, \dots, t} g(\mathbf{x}_k) \leq \frac{C_t}{\delta \sqrt{N_{t+1}}} = \mathcal{O}\left(\frac{1}{\delta \sqrt{t}}\right), \quad (15)$$

where $C_t = \max\{2h_0, L_t^{\max} \text{diam}(\mathcal{A})^2\}$ and $g = g^{\text{FW}}$ is the FW gap for AdaFW, AdaAFW, AdaPFW and $C_t = \text{radius}(\mathcal{A}) \sqrt{2h_0 L_{t+1}}$ and $g = g^{\text{MP}}$ is the MP gap for AdaMP.

5.2 Convex objectives

For convex objectives we will be able to improve the results of Theorem 1. We will first state the convergence results for FW variants and then for MP.

For adaptive FW variants, we will be able to give an $\mathcal{O}(1/\delta^2 t)$ convergence rate on the primal-dual gap, which trivially implies a bound on the objective suboptimality. In order to define the primal-dual gap, we

define the following dual objective function

$$\psi(\mathbf{u}) \stackrel{\text{def}}{=} -f^*(\mathbf{u}) - \sigma_{\mathcal{D}}(-\mathbf{u}), \quad (16)$$

where f^* denotes the convex conjugate of f and $\sigma_{\mathcal{D}}(\mathbf{x}) \stackrel{\text{def}}{=} \sup\{\mathbf{x} \cdot \mathbf{a} : \mathbf{a} \in \mathcal{D}\}$ is the support function over \mathcal{D} , which is the convex conjugate of the indicator function. Note that ψ is concave and that when f convex, we have by duality $\min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x}_t) = \max_{\mathbf{u} \in \mathbb{R}^p} \psi(\mathbf{u})$.

Theorem 2 (FW variants). *Let f be convex, \mathbf{x}_t denote the iterate generated by any of the proposed FW variants (AdaFW, AdaAFW, AdaPFW) after t iterations, with $N_t \geq 1$, and let \mathbf{u}_t be defined recursively as $\mathbf{u}_0 = \nabla f(\mathbf{x}_0)$, $\mathbf{u}_{t+1} = (1 - \xi_t)\mathbf{u}_t + \xi_t \nabla f(\mathbf{x}_t)$, where $\xi_t = 2/(\delta N_t + 2)$ if t is a good step and $\xi_t = 0$ otherwise. Then we have:*

$$h_t \leq f(\mathbf{x}_t) - \psi(\mathbf{u}_t) \quad (17)$$

$$\begin{aligned} &\leq \frac{2\bar{L}_t \text{diam}(\mathcal{A})^2}{\delta^2 N_t + \delta} + \frac{2(1 - \delta)}{\delta^2 N_t^2 + \delta N_t} (f(\mathbf{x}_0) - \psi(\mathbf{u}_0)) \\ &= \mathcal{O}\left(\frac{1}{\delta^2 t}\right). \end{aligned} \quad (18)$$

We will now give a similar sublinear convergence for Adaptive MP. As the similar results of (Locatello et al., 2018), it relies on the definition of *atomic norm* or gauge function of a set \mathcal{B} , defined as $\|\cdot\|_{\mathcal{B}} \stackrel{\text{def}}{=} \inf\{c > 0 : \mathbf{x} \in c \cdot \text{conv}(\mathcal{B})\}$.

Theorem 3 (MP). *Let f be convex, \mathbf{x}^* be an arbitrary solution to (OPT-MP) and let $R_{\mathcal{B}}$ the level set radius:*

$$R_{\mathcal{B}} = \max_{\substack{\mathbf{x} \in \text{lin}(\mathcal{A}) \\ f(\mathbf{x}) \leq f(\mathbf{x}_0)}} \|\mathbf{x} - \mathbf{x}^*\|_{\mathcal{B}}. \quad (19)$$

If we denote by \mathbf{x}_t the iterate generated by AdaMP after $t \geq 1$ iterations and $\beta = \delta/R_{\mathcal{B}}$, then we have:

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{2\bar{L}_t \text{radius}(\mathcal{A})^2}{\beta^2 t + \beta} + \frac{2(1 - \beta)}{\beta^2 t^2 + \beta t} h_0 \quad (20)$$

$$= \mathcal{O}\left(\frac{1}{\beta^2 t}\right). \quad (21)$$

5.3 Strongly convex objectives

The next result states the linear convergence of some algorithm variants and uses the notions of pyramidal width (PWidth) and minimal directional width (mDW) that have been developed in (Lacoste-Julien, 2016) and (Locatello et al., 2017) respectively, which we state in Appendix A for completeness. We note that the pyramidal width of a set \mathcal{A} is lower bounded by the minimal width over all subsets of atoms, and thus is strictly greater than zero if the number of atoms is

finite. The minimal directional width is a much simpler quantity and always strictly greater than zero by the symmetry of our domain.

Theorem 4 (Linear convergence rate for strongly convex objectives). *Let f be μ -strongly convex. Then for [AdaAFW](#), [AdaPFW](#) or [AdaMP](#) we have the following linear decrease for each good step t :*

$$h_{t+1} \leq (1 - \delta^2 \rho_t) h_t, \quad (22)$$

where

$$\rho_t = \frac{\mu}{4L_t} \left(\frac{\text{PWidth}(\mathcal{A})}{\text{diam}(\mathcal{A})} \right)^2 \text{ for } \text{AdaAFW} \text{ and } \text{AdaPFW},$$

$$\rho_t = \frac{\mu}{L_t} \left(\frac{\text{mDW}(\mathcal{A})}{\text{radius}(\mathcal{A})} \right)^2 \text{ for } \text{AdaMP}.$$

The previous theorem gives a geometric decrease on good steps. Combining this theorem with the bound for the number of bad steps in (9), and noting that the sufficient decrease guarantees that the objective is monotonically decreasing, we obtain a global linear convergence for [AdaAFW](#), [AdaPFW](#) and [AdaMP](#).

5.4 Discussion

Non-convex objectives. [Lacoste-Julien \(2016\)](#) studied the convergence of FW assuming the linear subproblems are solved exactly ($\delta = 1$) and obtained a rate of the form (14) with $C_0 = \max\{2h_0, L \text{diam}(\mathcal{D})^2\}$ instead. Both rates are similar, although our analysis is more general as it allows to consider the case in which linear subproblems are solved approximately ($\delta < 1$) and also gives rates for the Away-steps and Pairwise variants, for which no rates were previously known.

Theorem 1 also gives the first known convergence rates for a variant of MP on general non-convex functions. Contrary to the case of FW, this bound depends on the mean instead of the maximum of the Lipschitz estimate.

Convex objectives. Compared with ([Jaggi, 2013](#)), the primal-dual rates of Theorem 2 are stronger as they hold for the last iterate and not only for the minimum over previous iterates. To the best of our knowledge, primal-dual convergence rates on the last iterate have only been derived in ([Nesterov, 2017](#)) and were not extended to approximate linear subproblems nor the Away-steps and Pairwise variants.

Compared to [Nesterov \(2017\)](#) on the special case of exact subproblems ($\delta = 1$), the rates of Theorem 2 are similar but with \bar{L}_t replaced by L . Hence, in the regime $\bar{L}_t \ll L$ (as is often verified in practice), our bounds have a much smaller leading constant.

For MP, [Locatello et al. \(2018\)](#) obtains a similar convergence rate of the form $\mathcal{O}(L_B R_B^2 / (\delta^2 t))$, where L_B is

the Lipschitz constant of ∇f under the atomic norm, instead of the adaptive, averaged Lipschitz estimate in our case.

Strongly convex objectives. For the FW variants, the rates are identical to the ones in ([Lacoste-Julien and Jaggi, 2015](#), Theorem 1), but where L is replaced with the adaptive L_t in the linear rate factor, giving a larger per-iteration decrease whenever $L_t < L$. Our rates are the first also covering approximate subproblems for Away-Steps and Pairwise FW algorithms. It is also worth noticing that both Away-steps FW and Pairwise FW have only been previously analyzed in the presence of exact line search ([Lacoste-Julien and Jaggi, 2015](#)). Additionally, unlike ([Lacoste-Julien and Jaggi, 2015](#)), we do not require a smoothness assumption on f outside of the domain. Finally, for the case of MP, we again obtain the same convergence rates as in ([Locatello et al., 2017](#), Theorem 7), but with L replaced by L_t .

6 Experiments

We compared the proposed methods across three problems and three datasets. We show the main properties of these datasets in the table of Figure 1, where density denotes the fraction of nonzero coefficients in data matrix and where the last two columns are quantities that arise during the optimization of [AdaPFW](#) and shed light into their empirical value. In both cases t is the number of iterates until 10^{-10} suboptimality is achieved.

The first problem that we consider is a logistic regression with the ℓ_1 norm constraint on the coefficients $\|\mathbf{x}\|_1 \leq \beta$, where β is chosen to give approximately 1%, 20% of nonzero coefficients respectively. We applied this problem on two different datasets: [Madelon](#) and [RCV1](#) and show the results in columns Figure 1, subplots A, B, C, D. In that figure we also show the performance of FW, Away-steps FW (AFW) and Pairwise FW (PFW), all of them using the step-size $\gamma_t = \min\{g_t L^{-1} \|\mathbf{d}_t\|^{-2}, \gamma_t^{\max}\}$, as well as the adaptive step-size FW variants of [Dunn \(1980\)](#) and ([Beck et al., 2015](#)), which we denote D-FW and B-FW respectively.

The second problem that we consider is collaborative filtering. We used the [MovieLens 1M](#) dataset, which contains 1 million movie ratings, and consider the problem of minimizing a Huber loss, as in ([Mehta et al., 2007](#)), between the true ratings and a matrix \mathbf{X} . We also constrain the matrix by its nuclear norm $\|\mathbf{X}\|_* \leq \beta$, where β is chosen to give approximately 1% and 20% if non-zero singular values respectively. In this case the AFW and PFW variants were not considered as they are not directly applicable to this problem with as the size of the active set is potentially unbounded. The results of this comparison can be seen in subplots E

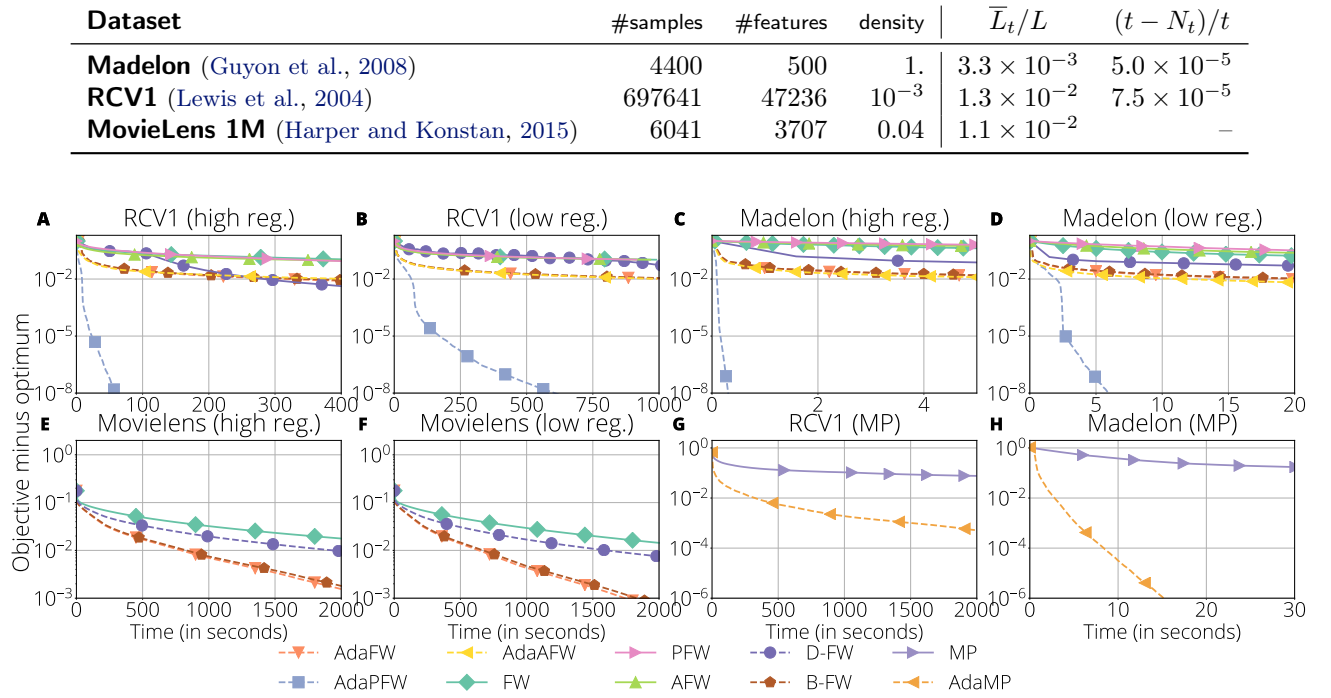


Figure 1: **Top table:** description of the datasets. **Bottom figure:** Benchmark of different FW and MP variants. Adaptive variants proposed in this paper are in dashed lines. Problem in A, B, C, D = logistic regression with ℓ_1 -constrained coefficients, in E, F = Huber regression with on the nuclear norm constrained coefficients and in G, H = unconstrained logistic regression (MP variants). In all the considered datasets and regularization regimes adaptive variants have a much faster convergence than non-adaptive ones.

and F of Figure 1.

The third problem that we consider is an (unconstrained) logistic regression problem that we solve using matching pursuit. In this case our atoms are the euclidean basis and so there is no explicit regularization. This is a common setting for MP, where the regularization comes from performing early stopping. Subplots G and H show the comparison between MP and AdaMP on the RCV1 and Madelon dataset. In all cases the linear subproblems were solved exactly ($\delta = 1$, machine precision in the case of the nuclear norm constrained problem).

We comment on a couple of observed trends from these results:

- **Adaptive vs non-adaptive.** Across the different datasets, problems and regularization regimes we found that adaptive step-size methods always perform better than their non-adaptive variant.
- **Pairwise FW.** AdaPFW shows a surprisingly good performance when it is applicable, specially in the high regularization regime. A possible interpretation for this is that it is the only variant of FW in which the coefficients associated with previous atoms are

not shrunk when adding a new atom, hence large step sizes are potentially even more beneficial as coefficients that are already close to optimal do not get necessarily modified in subsequent updates.

7 Conclusion and future work

In this work we have proposed and analyzed a novel adaptive step-size scheme that can be used in projection-free methods such as FW and MP. The method has minimal computational overhead and does not rely on any step-size hyperparameter (except for an initial estimate). Numerical experiments show large computational gains on a variety of problems.

A possible extension of this work is to develop adaptive step-size strategies for randomized variants of FW such as (Lacoste-Julien et al., 2013; Kerdreux et al., 2018; Mokhtari et al., 2018), in which there is stochasticity in the linear subproblems.

Another area of future research is to improve the convergence rate of the (adaptive) Pairwise FW method. Due to the very pessimistic bound on its number of bad steps, there is still a large gap between its excellent empirical performance and its known convergence rate.

Acknowledgements

The authors would like to thank Vlad Niculae for valuable feedback on the manuscript.

FP is funded through the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement 748900.

References

- Armijo, L. (1966). [Minimization of functions having Lipschitz continuous first partial derivatives](#). *Pacific Journal of Mathematics*.
- Beck, A., Pauwels, E., and Sabach, S. (2015). [The cyclic block conditional gradient method for convex optimization problems](#). *SIAM Journal on Optimization*.
- Beck, A. and Teboulle, M. (2009). [Gradient-based algorithms with applications to signal recovery](#). *Convex optimization in signal processing and communications*.
- Bertsekas, D. P. (1999). *Nonlinear programming*. Athena Scientific.
- Demyanov, V. and Rubinov, A. (1967). [The minimization of a smooth convex functional on a convex set](#). *SIAM Journal on Control*.
- Dunn, J. C. (1980). [Convergence rates for conditional gradient sequences generated by implicit step length rules](#). *SIAM Journal on Control and Optimization*.
- Frank, M. and Wolfe, P. (1956). [An algorithm for quadratic programming](#). *Naval Research Logistics (NRL)*.
- Garber, D. and Hazan, E. (2013). [A linearly convergent conditional gradient algorithm with applications to online and stochastic optimization](#). *arXiv preprint arXiv:1301.4666*.
- Goldstein, A. A. (1965). [On steepest descent](#). *Journal of the Society for Industrial and Applied Mathematics, Series A: Control*.
- Guélat, J. and Marcotte, P. (1986). [Some comments on Wolfe’s ‘away step’](#). *Mathematical Programming*.
- Guyon, I., Gunn, S., Nikravesh, M., and Zadeh, L. A. (2008). *Feature extraction: foundations and applications*, volume 207. Springer.
- Harper, F. M. and Konstan, J. A. (2015). [The movielens datasets: History and context](#). *ACM Transactions on Interactive Intelligent Systems (TiiS)*.
- Jaggi, M. (2013). [Revisiting Frank-Wolfe: projection-free sparse convex optimization](#). In *International Conference on Machine Learning*.
- Kerdreux, T., Pedregosa, F., and d’Aspremont, A. (2018). [Frank-Wolfe with Subsampling Oracle](#). In *Proceedings of the 35th International Conference on Machine Learning*, Proceedings of Machine Learning Research. PMLR.
- Lacoste-Julien, S. (2016). [Convergence rate of Frank-Wolfe for non-convex objectives](#). *arXiv preprint arXiv:1607.00345*.
- Lacoste-Julien, S. and Jaggi, M. (2015). [On the global linear convergence of Frank-Wolfe optimization variants](#). In *Advances in Neural Information Processing Systems*.
- Lacoste-Julien, S., Jaggi, M., Schmidt, M., and Pletscher, P. (2013). [Block-Coordinate Frank-Wolfe Optimization for Structural SVMs](#). In *Proceedings of the 30th International Conference on Machine Learning*.
- Lewis, D. D., Yang, Y., Rose, T. G., and Li, F. (2004). [RCV1: A new benchmark collection for text categorization research](#). *Journal of machine learning research*, 5(Apr):361–397.
- Locatello, F., Khanna, R., Tschannen, M., and Jaggi, M. (2017). [A Unified Optimization View on Generalized Matching Pursuit and Frank-Wolfe](#). In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*.
- Locatello, F., Raj, A., Karimireddy, S. P., Raetsch, G., Schölkopf, B., Stich, S., and Jaggi, M. (2018). [On Matching Pursuit and Coordinate Descent](#). In *Proceedings of the 35th International Conference on Machine Learning*.
- Mallat, S. G. and Zhang, Z. (1993). [Matching pursuits with time-frequency dictionaries](#). *IEEE Transactions on signal processing*.
- Mehta, B., Hofmann, T., and Nejd, W. (2007). [Robust collaborative filtering](#). In *Proceedings of the 2007 ACM conference on Recommender systems*. ACM.
- Mitchell, B., Demyanov, V. F., and Malozemov, V. (1974). [Finding the point of a polyhedron closest to the origin](#). *SIAM Journal on Control*.
- Mokhtari, A., Hassani, H., and Karbasi, A. (2018). [Stochastic Conditional Gradient Methods: From Convex Minimization to Submodular Maximization](#). *arXiv*.
- Nesterov, Y. (2013). [Gradient methods for minimizing composite functions](#). *Mathematical Programming*.
- Nesterov, Y. (2017). [Complexity bounds for primal-dual methods minimizing the model of objective function](#). *Mathematical Programming*.
- Reddi, S. J., Sra, S., Póczos, B., and Smola, A. (2016). [Stochastic Frank-Wolfe methods for nonconvex op-](#)

timization. In *54th Annual Allerton Conference on Communication, Control, and Computing*.

Wolfe, P. (1969). Convergence conditions for ascent methods. *SIAM review*.

Step-Size Adaptivity in Projection-Free Optimization

Supplementary material

Outline. The supplementary material of this paper is organized as follows.

- [Appendix A](#) contains definitions and properties relative to the objective function and/or the domain, such as the definition of geometric strong convexity and pyramidal width.
- [Appendix B](#) we present key inequalities on the abstract algorithm which are used by the different convergence proofs.
- [Appendix C](#) provides a proof of convergence for non-convex objectives (Theorem 1).
- [Appendix D](#) provides a proof of convergence for convex objectives (Theorem 2).
- [Appendix E](#) provides a proof of linear convergence for all variants except FW (Theorem 4).

Appendix A Basic definitions and properties

In this section we give basic definitions and properties relative to the objective function and/or the domain, such as the definition of geometric strong convexity and pyramidal width. These definitions are not specific to our algorithms and have appeared in different sources such as [Lacoste-Julien and Jaggi \(2015\)](#); [Locatello et al. \(2017\)](#). We merely gather them here for completeness.

Definition 1 (Geometric strong convexity). *We define the **geometric strong convexity constant** μ_f^A as*

$$\mu_f^A \stackrel{\text{def}}{=} \inf_{\substack{\mathbf{x}, \mathbf{x}^* \in \mathcal{D} \\ \langle \nabla f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle < 0}} \frac{2}{\gamma(\mathbf{x}, \mathbf{x}^*)^2} \left(f(\mathbf{x}^*) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle \right) \quad (23)$$

$$\text{where } \gamma(\mathbf{x}, \mathbf{x}^*) \stackrel{\text{def}}{=} \frac{\langle -\nabla f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle}{\langle -\nabla f(\mathbf{x}), \mathbf{s}_f(\mathbf{x}) - \mathbf{v}_f(\mathbf{x}) \rangle}, \quad (24)$$

where

$$\mathbf{s}_f(\mathbf{x}) \stackrel{\text{def}}{=} \arg \min_{\mathbf{v} \in \mathcal{A}} \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle \quad (25)$$

$$\mathbf{v}_f(\mathbf{x}) \stackrel{\text{def}}{=} \arg \min_{\substack{\mathbf{v} = \mathbf{v}_S(\mathbf{x}) \\ S \in \mathcal{S}_\mathbf{x}}} \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle \quad (26)$$

$$\mathbf{v}_S(\mathbf{x}) \stackrel{\text{def}}{=} \arg \max_{\mathbf{v} \in S} \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle \quad (27)$$

where $\mathcal{S} \subseteq \mathcal{A}$ and $\mathcal{S}_\mathbf{x} \stackrel{\text{def}}{=} \{S | S \subseteq \mathcal{A} \text{ such that } \mathbf{x} \text{ is a proper convex combination of all the elements in } S\}$ (recall \mathbf{x} is a proper convex combination of elements in S when $\mathbf{x} = \sum_i \alpha_i \mathbf{s}_i$ where $\mathbf{s}_i \in S$ and $\alpha_i \in (0, 1)$).

Definition 2 (Pyramidal width). *The **pyramidal width** of a set \mathcal{A} is the smallest pyramidal width of all its faces, i.e.*

$$\text{PWidth}(\mathcal{A}) \stackrel{\text{def}}{=} \min_{\substack{\mathbf{x} \in \mathcal{K} \\ \mathcal{K} \in \text{faces}(\text{conv}(\mathcal{A})) \\ \mathbf{r} \in \text{cone}(\mathcal{K} - \mathbf{x}) \setminus \{0\}}} \text{PdirW}(\mathcal{K} \cap \mathcal{A}, \mathbf{r}, \mathbf{x}) \quad (28)$$

where PdirW is the pyramidal directional width, defined as

$$\text{PdirW}(W)(\mathcal{A}, \mathbf{r}, \mathbf{x}) \stackrel{\text{def}}{=} \min_{S \in \mathcal{S}_\mathbf{x}} \max_{\mathbf{s} \in \mathcal{A}, \mathbf{v} \in S} \left\langle \frac{\mathbf{r}}{\|\mathbf{r}\|_2}, \mathbf{s} - \mathbf{v} \right\rangle \quad (29)$$

We now relate these two geometric quantities together.

Lemma 1 (Lower bounding μ_f^A). *Let f μ -strongly convex on $\mathcal{D} = \text{conv}(\mathcal{A})$. Then*

$$\mu_f^A \geq \mu \cdot (\text{PWidth}(\mathcal{A}))^2 \quad (30)$$

Proof. We refer to (Lacoste-Julien and Jaggi, 2015, Theorem 6). \square

Proposition 1. $\text{PWidth}(\mathcal{A}) \leq \text{diam}(\text{conv}(\mathcal{A}))$ where $\text{diam}(\mathcal{X}) \stackrel{\text{def}}{=} \sup_{x,y \in \mathcal{X}} \|x - y\|_2$.

Proof. First note that given $\mathbf{r} \in \mathcal{R}$, $\mathbf{s} \in \mathcal{S}$, $\mathbf{v} \in \mathcal{V}$ with $\mathcal{R}, \mathcal{S}, \mathcal{V} \subseteq \mathbb{R}^n$, we have

$$\langle \mathbf{r} / \|\mathbf{r}\|_2, \mathbf{s} - \mathbf{v} \rangle \leq \|\mathbf{s} - \mathbf{v}\|_2 \quad \forall \mathbf{r} \in \mathcal{R}, \mathbf{s} \in \mathcal{S}, \mathbf{v} \in \mathcal{V} \quad (31)$$

$$\Rightarrow \max_{\mathbf{s} \in \mathcal{S}, \mathbf{v} \in \mathcal{V}} \langle \mathbf{r} / \|\mathbf{r}\|_2, \mathbf{s} - \mathbf{v} \rangle \leq \max_{\mathbf{s} \in \mathcal{S}, \mathbf{v} \in \mathcal{V}} \|\mathbf{s} - \mathbf{v}\|_2 \quad \forall \mathbf{r} \in \mathcal{R} \quad (32)$$

$$\Rightarrow \min_{\mathbf{r} \in \mathcal{R}} \max_{\mathbf{s} \in \mathcal{S}, \mathbf{v} \in \mathcal{V}} \langle \mathbf{r} / \|\mathbf{r}\|_2, \mathbf{s} - \mathbf{v} \rangle \leq \max_{\mathbf{s} \in \mathcal{S}, \mathbf{v} \in \mathcal{V}} \|\mathbf{s} - \mathbf{v}\|_2 \quad (33)$$

Applying this result to the definition of pyramidal width we have

$$\text{PWidth}(\mathcal{A}) = \min_{\substack{\mathbf{x} \in \mathcal{K} \\ \mathcal{K} \in \text{faces}(\text{conv}(\mathcal{A})) \\ \mathbf{r} \in \text{cone}(\mathcal{K} - \mathbf{x}) \setminus \{0\}}} \text{PdirW}(\mathcal{K} \cap \mathcal{A}, \mathbf{r}, \mathbf{x}) \quad (34)$$

$$= \min_{\substack{\mathbf{x} \in \mathcal{K} \\ \mathcal{K} \in \text{faces}(\text{conv}(\mathcal{A})) \\ \mathbf{r} \in \text{cone}(\mathcal{K} - \mathbf{x}) \setminus \{0\}}} \min_{\mathcal{S} \subseteq \mathcal{S}_x} \max_{\mathbf{s} \in \mathcal{A}, \mathbf{v} \in \mathcal{S}} \left\langle \frac{\mathbf{r}}{\|\mathbf{r}\|}, \mathbf{s} - \mathbf{v} \right\rangle \quad (35)$$

$$= \min_{\mathbf{r} \in \mathcal{R}} \max_{\mathbf{s} \in \mathcal{A}, \mathbf{v} \in \mathcal{V}} \left\langle \frac{\mathbf{r}}{\|\mathbf{r}\|}, \mathbf{s} - \mathbf{v} \right\rangle \quad (36)$$

$$(37)$$

where $\mathcal{R} = \{\text{cone}(\mathcal{K} - \mathbf{x}) \setminus \{0\} : \text{for some } \mathbf{x} \in \mathcal{K}, \mathcal{K} \in \text{faces}(\text{conv}(\mathcal{A}))\}$ and \mathcal{V} is some subset of \mathcal{A} . Applying the derived result we have that

$$\begin{aligned} \text{PWidth}(\mathcal{A}) &\leq \max_{\mathbf{s} \in \mathcal{A}, \mathbf{v} \in \mathcal{V}} \|\mathbf{s} - \mathbf{v}\|_2 \\ &\leq \max_{\mathbf{s}, \mathbf{v} \in \text{conv}(\mathcal{A})} \|\mathbf{s} - \mathbf{v}\|_2 \\ &= \text{diam}(\text{conv}(\mathcal{A})) \end{aligned}$$

\square

Definition 3. The *minimal directional width* $\text{mDW}(\mathcal{A})$ of a set of atoms \mathcal{A} is defined as

$$\text{mDW}(\mathcal{A}) = \min_{\mathbf{d} \in \text{lin}(\mathcal{A})} \max_{\mathbf{z} \in \mathcal{A}} \frac{\langle \mathbf{z}, \mathbf{d} \rangle}{\|\mathbf{d}\|}. \quad (38)$$

Note that in contrast to the pyramidal width, the minimal directional width here is a much simpler and robust property of the atom set \mathcal{A} , not depending on its combinatorial face structure of the polytope. As can be seen directly from the definition above, the $\text{mDW}(\mathcal{A})$ is robust when adding a duplicate atom or small perturbation of it to \mathcal{A} .

Appendix B Preliminaries: Key Inequalities

In this appendix we prove that the sufficient decrease condition verifies a recursive inequality. This key result is used by all convergence proofs.

Lemma 2. *The following inequality is verified for all proposed algorithms (with $\gamma_t^{\max} = +\infty$ for [AdaMP](#)):*

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \xi g_t + \frac{\xi^2 L_t}{2} \|\mathbf{d}_t\|^2 \quad \text{for all } \xi \in [0, \gamma_t^{\max}]. \quad (39)$$

Proof. We start the proof by proving an optimality condition of the step-size. Consider the following quadratic optimization problem:

$$\underset{\xi \in [0, \gamma_t^{\max}]}{\text{minimize}} \quad -\xi g_t + \frac{L_t \xi^2}{2} \|\mathbf{d}_t\|^2. \quad (40)$$

Deriving with respect to ξ and noting that on all the considered algorithms we have $\langle \nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle \leq 0$, one can easily verify that the global minimizer is achieved at the value

$$\min \left\{ \frac{g_t}{L_t \|\mathbf{d}_t\|^2}, \gamma_t^{\max} \right\}, \quad (41)$$

where $g_t = \langle -\nabla f(\mathbf{x}), \mathbf{d}_t \rangle$. This coincides with the value of γ_{t+1} computed by the backtracking procedure on the different algorithms and so we have:

$$-\gamma_t g_t + \frac{L_t \gamma_t^2}{2} \|\mathbf{d}_t\|^2 \leq -\xi g_t + \frac{L_t \xi^2}{2} \|\mathbf{d}_t\|^2 \quad \text{for all } \xi \in [0, \gamma_t^{\max}]. \quad (42)$$

We can now write the following sequence of inequalities, that combines the sufficient decrease condition with this last inequality:

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \gamma_t g_t + \frac{L_t \gamma_t^2}{2} \|\mathbf{d}_t\|^2 \quad (43)$$

$$\stackrel{(40)}{\leq} f(\mathbf{x}_t) - \xi g_t + \frac{L_t \xi^2}{2} \|\mathbf{d}_t\|^2 \quad \text{for any } \xi \in [0, \gamma_t^{\max}]. \quad (44)$$

□

Proposition 2. *The Lipschitz estimate L_t is bounded as $L_t \leq \max\{\tau L, L_{-1}\}$.*

Proof. If the sufficient decrease condition is verified then we have $L_t = L_{t-1}/\eta$ and so $L_t \leq L_{t-1}$. If its not, then we must have $L_{t-1}/\eta \leq L$, and the Lipschitz estimate cannot larger than τL . Combining both bounds we obtain

$$L_t \leq \max\{\tau L, L_{t-1}\}. \quad (45)$$

Applying the same bound recursively on L_{t-1} leads to the claimed bound $L_t \leq \max\{\tau L, L_{-1}/\eta\}$. □

Lemma 3. *Let $g(\cdot)$ be as in Theorem 1, i.e., $g(\cdot) = g^{FW}(\cdot)$ for FW variants ([AdaFW](#), [AdaAFW](#), [AdaPFW](#)) and $g(\cdot) = g^{MP}(\cdot)$ for MP variants ([AdaMP](#)). Then for any of these algorithms we have*

$$g_t \geq \delta g(\mathbf{x}_t). \quad (46)$$

Proof. • For [AdaFW](#) and [AdaMP](#), Eq. (46) follows immediately from the definition of g_t and $g(\mathbf{x}_t)$.

• For [AdaAFW](#), by the way the descent direction is selected in Line 6, we always have

$$g_t \geq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{s}_t \rangle \geq \delta g(\mathbf{x}_t), \quad (47)$$

where the last inequality follows from the definition of \mathbf{s}_t

- For [AdaPFW](#), we have

$$g_t = \langle \nabla f(\mathbf{x}_t), \mathbf{v}_t - \mathbf{s}_t \rangle = \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{s}_t \rangle + \langle \nabla f(\mathbf{x}_t), \mathbf{v}_t - \mathbf{x}_t \rangle \quad (48)$$

$$\geq \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{s}_t \rangle \geq \delta g(\mathbf{x}_t) \quad (49)$$

where the term $\langle \nabla f(\mathbf{x}_t), \mathbf{v}_t - \mathbf{x}_t \rangle$ is positive by definition of \mathbf{v}_t since \mathbf{x}_t is necessarily in the convex envelope of \mathcal{S}_t . The second inequality follows from the definition of \mathbf{s}_t . \square

Lemma 4. *Let N_t be the total number of evaluations of the sufficient decrease condition up to iteration t . Then we have*

$$N_t \leq \sum_{i=0}^t n_i = \left[1 + \frac{\log \eta}{\log \tau} \right] (t+1) + \frac{1}{\log \tau} \max \left\{ \log \frac{\tau L}{L_{-1}}, 0 \right\}, \quad (50)$$

where $(a)_+ = \max\{a, 0\}$

Proof. This proof follows roughly that of ([Nesterov, 2013](#), Lemma 3), albeit with a slightly different bound on L_t due to algorithmic differences.

Denote by $n_i \geq 1$ the number of evaluations of the sufficient decrease condition. Since the algorithm multiplies by τ every time that the sufficient condition is not verified, we have

$$L_i = \frac{1}{\eta} L_{i-1} \tau^{n_i-1}. \quad (51)$$

Taking logarithms on both sides we obtain

$$n_i \leq 1 + \frac{\log \eta}{\log \tau} + \frac{1}{\tau} \log \frac{L_i}{L_{i-1}}. \quad (52)$$

Summing from $i = 0$ to $i = t$ gives

$$N_t \leq \sum_{i=0}^t n_i = \left[1 + \frac{\log \eta}{\log \tau} \right] (t+1) + \frac{1}{\log \tau} \log \left(\frac{L_t}{L_{-1}} \right) \quad (53)$$

Finally, from Proposition 2 we have the bound $L_t \leq \max\{\tau L, L_{-1}\}$, which we can use to bound the numerator's last term. This gives the claimed bound

$$N_t \leq \sum_{i=0}^t n_i = \left[1 + \frac{\log \eta}{\log \tau} \right] (t+1) + \frac{1}{\log \tau} \max \left\{ \log \frac{\tau L}{L_{-1}}, 0 \right\}. \quad (54)$$

\square

Appendix B.1 A bound on the number of bad steps

To prove the linear rates for the adaptive AFW and adaptive PFW algorithm it is necessary to bound the number of bad steps. There are two different types of bad steps: “drop” steps and “swap” steps. These names come from how the active set \mathcal{S}_t changes. In a drop step, an atom is removed from the active set (i.e. $|\mathcal{S}_{t+1}| < |\mathcal{S}_t|$). In a swap step, the size of the active set remains unchanged (i.e. $|\mathcal{S}_{t+1}| = |\mathcal{S}_t|$) but one atom is swapped with another one not in the active set. Note that drop steps can occur in the (adaptive) Away-steps and Pairwise, but swap steps can only occur in the Pairwise variant.

For the proofs of linear convergence in [Appendix E](#), we show that these two types of bad steps are only problematic when $\gamma_t = \gamma_t^{\max} < 1$. In these scenarios, we cannot provide a meaningful decrease bound. However, we show that the number of bad steps we take is bounded. The following two lemmas adopted from ([Lacoste-Julien and Jaggi, 2015](#), Appendix C) bound the number of drop steps and swap steps the adaptive algorithms can take.

Lemma 5. *After T steps of [AdaAFW](#) or [AdaPFW](#), there can only be $T/2$ drop steps. Also, if there is a drop step at step $t+1$, then $f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) < 0$.*

Proof. Let A_t denote the number of steps that added a vertex in the expansion, and let D_t be the number of drop steps. Then $1 \leq |\mathcal{S}_t| = |\mathcal{S}_0| + A_t - D_t$ and we clearly have $A_t - D_t \leq t$. Combining these two inequalities we have that $D_t \leq \frac{1}{2}(|\mathcal{S}_0| - 1 + t) = \frac{t}{2}$.

To show $f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) < 0$, because of Lemma 2, it suffices to show that

$$-\gamma_t g_t + \frac{1}{2} \gamma_t^2 L_t \|\mathbf{d}_t\|^2 < 0, \quad (55)$$

with $\gamma_t = \gamma_t^{\max}$ (recall drop steps only occur when $\gamma_t = \gamma_t^{\max}$). Note this is a convex quadratic in γ_t which is precisely less than or equal to 0 when $\gamma_t \in [0, 2g_t/L_t \|\mathbf{d}_t\|^2]$. Thus in order to show $f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) < 0$ it suffices to show $\gamma_t^{\max} \in (0, 2g_t/L_t \|\mathbf{d}_t\|^2]$. This follows immediately since $0 < \gamma_t^{\max} \leq g_t/L_t \|\mathbf{d}_t\|^2$. \square

Since in the AdaAFW algorithm all bad steps are drop steps, the previous lemma implies that we can effectively bound the number of bad steps by $t/2$, which is the bound claimed in (9).

Lemma 6. *There are at most $3|\mathcal{A}|!$ bad steps between any two good steps in AdaPFW. Also, if there is a swap step at step $t + 1$, then $f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) < 0$.*

Proof. Note that bad steps only occur when $\gamma_t = \gamma_t^{\max} = \alpha_{\mathbf{v}_t, t}$. When this happens there are two possibilities; we either move all the mass from \mathbf{v}_t to a new atom $\mathbf{s}_t \notin \mathcal{S}_t$ (i.e. $\alpha_{\mathbf{v}_t, t+1} = 0$ and $\alpha_{\mathbf{s}_t, t+1} = \alpha_{\mathbf{v}_t, t}$) and preserve the cardinality of our active set ($|\mathcal{S}_{t+1}| = |\mathcal{S}_t|$) or we move all the mass from \mathbf{v}_t to an old atom $\mathbf{s}_t \in \mathcal{S}_t$ (i.e. $\alpha_{\mathbf{s}_t, t+1} = \alpha_{\mathbf{s}_t, t} + \alpha_{\mathbf{v}_t, t}$) and the cardinality of our active set decreases by 1 ($|\mathcal{S}_{t+1}| < |\mathcal{S}_t|$). In the former case, the possible values of the coordinates $\alpha_{\mathbf{v}}$ do not change, but they are simply rearranged in the possible $|\mathcal{A}|$ slots. Note further every time the mass from \mathbf{v}_t moves to a new atom $\mathbf{s}_t \notin \mathcal{S}_t$ we have strict descent, i.e. $f(\mathbf{x}_{t+1}) < f(\mathbf{x}_t)$ unless \mathbf{x}_t is already optimal (see Lemma 5) and hence we cannot revisit the same point unless we have converged. Thus the maximum number of possible consecutive swap steps is bounded by the number of ways we can assign $|\mathcal{S}_t|$ numbers in $|\mathcal{A}|$ slots, which is $|\mathcal{A}|!/(|\mathcal{A}| - |\mathcal{S}_t|)!$. Furthermore, when the cardinality of our active set drops, in the worst case we will do a maximum number of drop steps before reducing the cardinality of our active set again. Thus starting with $|\mathcal{S}_t| = r$ the maximum number of bad steps B without making any good steps is upper bounded by

$$B \leq \sum_{k=1}^r \frac{|\mathcal{A}|!}{(|\mathcal{A}| - k)!} \leq |\mathcal{A}|! \sum_{k=0}^{\infty} \frac{1}{k!} = |\mathcal{A}|! e \leq 3|\mathcal{A}|!$$

\square

Appendix C Proofs of convergence for non-convex objectives

In this appendix we provide the convergence proof of Theorem 1. Although this theorem provides a unified convergence proof for both variants of FW and MP, for convenience we split the proof into one for FW variants (Theorem 1.A) and another one for variants of MP (Theorem 1.B)

Theorem 1.A. *Let \mathbf{x}_t denote the iterate generated by either [AdaFW](#), [AdaAFW](#) or [AdaPFW](#) after t iterations. Then for any iteration t with $N_{t+1} \geq 0$, we have the following suboptimality bound in terms of the FW gap:*

$$\lim_{k \rightarrow \infty} g^{FW}(\mathbf{x}_k) = 0 \quad \text{and} \quad \min_{k=0, \dots, t} g^{FW}(\mathbf{x}_k) \leq \frac{\max\{2h_0, L_t^{\max} \text{diam}(\mathcal{A})^2\}}{\delta \sqrt{N_{t+1}}} = \mathcal{O}\left(\frac{1}{\delta \sqrt{t}}\right) \quad (56)$$

Proof. By Lemma 2 we have the following inequality for any k and any $\xi \in [0, \gamma_k^{\max}]$,

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \xi g_k + \frac{\xi^2 C_k}{2}, \quad (57)$$

where we define $C_k \stackrel{\text{def}}{=} L_k \|\mathbf{d}_k\|^2$ for convenience. We consider now different cases according to the relative values of γ_k and γ_k^{\max} , yielding different upper bounds for the right hand side.

Case 1: $\gamma_k < \gamma_k^{\max}$

In this case, γ_k maximizes the right hand side of the (unconstrained) quadratic in inequality (57) which then becomes:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \frac{g_k^2}{2C_k} \leq f(\mathbf{x}_k) - \frac{g_k}{2} \min\left\{\frac{g_k}{C_k}, 1\right\} \quad (58)$$

Case 2: $\gamma_k = \gamma_k^{\max} \geq 1$

By the definition of γ_t , this case implies that $C_k \leq g_k$ and so using $\xi = 1$ in (57) gives

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leq -g_k + \frac{C_k}{2} \leq -\frac{g_k}{2}. \quad (59)$$

Case 3: $\gamma_k = \gamma_k^{\max} < 1$

This corresponds to the problematic drop steps for [AdaAFW](#) or possibly swap steps for [AdaPFW](#), in which we will only be able to guarantee that the iterates are non-increasing. Choosing $\xi = 0$ in (57) we can at least guarantee that the objective function is non-increasing:

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) < 0. \quad (60)$$

Combining the previous cases. We can combine the inequalities obtained for the previous cases into the following inequality, valid for all $k \leq t$,

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leq -\frac{g_k}{2} \min\left\{\frac{g_k}{C_k}, 1\right\} \mathbb{1}\{k \text{ is a good step}\} \quad (61)$$

Adding the previous inequality from $k = 0$ up to t and rearranging we obtain

$$f(\mathbf{x}_0) - f(\mathbf{x}_{t+1}) \geq \sum_{k=0}^t \frac{g_k}{2} \min\left\{\frac{g_k}{L_k \|\mathbf{d}_k\|^2}, 1\right\} \mathbb{1}\{k \text{ is a good step}\} \quad (62)$$

$$\geq \sum_{k=0}^t \frac{g_k}{2} \min\left\{\frac{g_k}{C_k^{\max}}, 1\right\} \mathbb{1}\{k \text{ is a good step}\} \quad (63)$$

with $C_t^{\max} \stackrel{\text{def}}{=} L_t^{\max} \text{diam}(\mathcal{D})^2$. Taking the limit for $t \rightarrow +\infty$ we obtain that the right hand side is bounded by the compactness assumption on the domain \mathcal{D} and L -smoothness on f . The left hand side is an infinite sum, and so

a necessary condition for it to be bounded is that $g_k \rightarrow 0$, since $g_k \geq 0$ for all k . We have hence proven that $\lim_{k \rightarrow \infty} g_k = 0$, which by Lemma 3 implies $\lim_{k \rightarrow \infty} g(\mathbf{x}_k) = 0$. This proves the first claim of the Theorem.

We will now aim to derive explicit convergence rates for convergence towards a stationary point. Let $\tilde{g}_t = \min_{0 \leq k \leq t} g_k$, then from Eq. (63) we have

$$f(\mathbf{x}_0) - f(\mathbf{x}_{t+1}) \geq \sum_{k=0}^t \frac{\tilde{g}_t}{2} \min \left\{ \frac{\tilde{g}_t}{C_t^{\max}}, 1 \right\} \mathbb{1}\{k \text{ is a good step}\} \quad (64)$$

$$= N_{t+1} \frac{\tilde{g}_t}{2} \min \left\{ \frac{\tilde{g}_t}{C_t^{\max}}, 1 \right\}. \quad (65)$$

We now make a distinction of cases for the quantities inside the min.

- If $\tilde{g}_t \leq C_t^{\max}$, then (65) gives $f(\mathbf{x}_0) - f(\mathbf{x}_{t+1}) \geq N_{t+1} \tilde{g}_t^2 / (2C_t^{\max})$, which reordering gives

$$\tilde{g}_t \leq \sqrt{\frac{2C_t^{\max}(f(\mathbf{x}_0) - f(\mathbf{x}_{t+1}))}{N_{t+1}}} \leq \sqrt{\frac{2C_t^{\max}h_0}{N_{t+1}}} \leq \frac{2h_0 + C_t^{\max}}{2\sqrt{N_{t+1}}} \leq \frac{\max\{2h_0, C_t^{\max}\}}{\sqrt{N_{t+1}}}. \quad (66)$$

where in the third inequality we have used the inequality $\sqrt{ab} \leq \frac{a+b}{2}$ with $a = \sqrt{2h_0}$, $b = \sqrt{C_t^{\max}}$.

- If $\tilde{g}_t > C_t^{\max}$ we can get a better $\frac{1}{N_t}$ rate, trivially bounded by $\frac{1}{\sqrt{N_t}}$.

$$\tilde{g}_t \leq \frac{2h_0}{N_{t+1}} \leq \frac{2h_0}{\sqrt{N_{t+1}}} \leq \frac{\max\{2h_0, C_t^{\max}\}}{\sqrt{N_{t+1}}}. \quad (67)$$

We have obtained the same bound in both cases, hence we always have

$$\tilde{g}_t \leq \frac{\max\{2h_0, C_t^{\max}\}}{\sqrt{N_{t+1}}}. \quad (68)$$

Finally, from Lemma 3 we have $g(\mathbf{x}_k) \leq \frac{1}{\delta} g_k$ for all k and so

$$\min_{0 \leq k \leq t} g(\mathbf{x}_k) \leq \frac{1}{\delta} \min_{0 \leq k \leq t} g_k = \frac{1}{\delta} \tilde{g}_t \leq \frac{\max\{2h_0, C_t^{\max}\}}{\delta \sqrt{N_{t+1}}}, \quad (69)$$

and the claimed bound follows by definition of C_t^{\max} . The $\mathcal{O}(1/\delta\sqrt{t})$ rate comes from the fact that both \bar{L}_t and h_0 are upper bounded. \bar{L}_t is bounded by Proposition 2 and h_0 is bounded by assumption. \square

Theorem 1.B. *Let \mathbf{x}_t denote the iterate generated by AdaMP after t iterations. Then for $t \geq 0$ we have the following suboptimality bound in terms of the MP gap:*

$$\lim_{k \rightarrow \infty} g^{MP}(\mathbf{x}_k) = 0 \quad \text{and} \quad \min_{0 \leq k \leq t} g^{MP}(\mathbf{x}_k) \leq \frac{\text{radius}(\mathcal{A})}{\delta} \sqrt{\frac{2h_0 \bar{L}_t}{t+1}} = \mathcal{O}\left(\frac{1}{\delta\sqrt{t}}\right). \quad (70)$$

Proof. The proof similar than that of Theorem 1.A, except that in this case the expression of the step-size is simpler and does not depend on the minimum of two quantities. This avoids the case distinction that was necessary in the previous proof, resulting in a much simpler proof.

For all $k = 0, \dots, t$, using the sufficient decrease condition, and the definitions of γ_k and g_k :

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leq \gamma_k \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle + \frac{\gamma_k^2 L_k}{2} \|\mathbf{d}_k\|^2 \quad (71)$$

$$\leq \min_{\eta \geq 0} \left\{ -\eta g_k + \frac{1}{2} \eta^2 L_k \|\mathbf{d}_k\|^2 \right\} \quad (72)$$

$$\leq -\frac{g_k^2}{2L_k \|\mathbf{d}_k\|^2}, \quad (73)$$

where the last inequality comes from minimizing with respect to η . Summation over k from 0 to t and negating the previous inequality, we obtain:

$$\sum_{0 \leq k \leq t} \frac{g_k^2}{L_k} \leq (f(\mathbf{x}_0) - f(\mathbf{x}_t)) \text{radius}(\mathcal{A})^2 \leq 2h_0 \text{radius}(\mathcal{A})^2. \quad (74)$$

Taking the limit for $t \rightarrow \infty$ we obtain that the left hand side has a finite sum since the right hand side is bounded by assumption. Therefore, $g_k \rightarrow 0$, which by Lemma 3 implies $\lim_{k \rightarrow \infty} g(\mathbf{x}_k) = 0$. This proves the first claim of the Theorem.

We now aim to derive explicit convergence rates. Taking the min over the g_k s and taking a square root for the last inequality

$$\min_{0 \leq k \leq t} g_k \leq \sqrt{\frac{2h_0 \text{radius}(\mathcal{A})^2}{\sum_{0 \leq k \leq t} L_k^{-1}}} \quad (75)$$

The term $\left(n / \sum_{0 \leq k \leq t} L_k^{-1}\right)$ is the *harmonic mean* of the L_k s, which is always upper bounded by the average \bar{L}_t . Hence we obtain

$$\min_{0 \leq k \leq t} g_k \leq \frac{\text{radius}(\mathcal{A})}{\delta} \sqrt{\frac{2h_0 \bar{L}_t}{t+1}}. \quad (76)$$

The claimed rate then follows from using the bound $g(\mathbf{x}_k) \leq \frac{1}{\delta} g_k$ from Lemma 3, valid for all $k \geq 0$.

The $\mathcal{O}(1/\delta\sqrt{t})$ rate comes from the fact that both \bar{L}_t and h_0 are upper bounded. \bar{L}_t is bounded by Proposition 2 and h_0 is bounded by assumption.

□

Note: Harmonic mean vs arithmetic mean. The convergence rate for MP on non-convex objectives (Theorem 1) also holds by replacing \bar{L}_t by its harmonic mean $H_t \stackrel{\text{def}}{=} N_t / (\sum_{k=0}^{t-1} L_k^{-1} \mathbf{1}\{k \text{ is a good step}\})$ respectively. The harmonic mean is always less than the arithmetic mean, i.e., $H_t \leq \bar{L}_t$, although for simplicity we only stated both theorems with the arithmetic mean. Note that the Harmonic mean is Schur-concave, implying that $H_t \leq t \min\{L_k : k \geq t\}$, i.e. it is controlled by the smallest Lipschitz estimate encountered so far.

Appendix D Proofs of convergence for convex objectives

In this section we provide a proof the convergence rates stated in the theorem for convex objectives (Theorem 2). The section is structured as follows. We start by proving a technical result which is a slight variation of Lemma 2 and which will be used in the proof of Theorem 2. This is followed by the proof of Theorem 2.

Appendix D.1 Frank-Wolfe variants

Lemma 7. *For any of the proposed FW variants, if t is a good step, then we have*

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \xi g_t + \frac{\xi^2 L_t}{2} \|\mathbf{d}_t\|^2 \quad \text{for all } \xi \in [0, 1]. \quad (77)$$

Proof. If $\gamma_t^{\max} \geq 1$, the result is obvious from Lemma 2. If $\gamma_t^{\max} < 1$, then the inequality is only valid in the smaller interval $[0, \gamma_t^{\max}]$. However, since we have assumed that this is a good step, if $\gamma_t^{\max} < 1$ then we must have $\gamma_t < \gamma_t^{\max}$. By Lemma 2, we have

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) + \min_{\xi \in [0, \gamma_t^{\max}]} \left\{ \xi \langle \nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle + \frac{L_t \xi^2}{2} \|\mathbf{d}_t\|^2 \right\} \quad (78)$$

Because $\gamma_t < \gamma_t^{\max}$ and since the expression inside the minimization term of the previous equation is a quadratic function of ξ , γ_t is the unconstrained minimum and so we have

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) + \min_{\xi \geq 0} \left\{ \xi \langle \nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle + \frac{L_t \xi^2}{2} \|\mathbf{d}_t\|^2 \right\} \quad (79)$$

$$\leq f(\mathbf{x}_t) + \min_{\xi \in [0, 1]} \left\{ \xi \langle \nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle + \frac{L_t \xi^2}{2} \|\mathbf{d}_t\|^2 \right\}. \quad (80)$$

The claimed bound then follows from the optimality of the min. \square

The following lemma allows to relate the quantity $\langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{s}_t \rangle$ with a primal-dual gap and will be essential in the proof of Theorem 2.

Lemma 8. *Let \mathbf{s}_t be as defined in any of the FW variants. Then for any iterate $t \geq 0$ we have*

$$\langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{s}_t \rangle \geq \delta(f(\mathbf{x}_t) - \psi(\nabla f(\mathbf{x}_t))). \quad (81)$$

Proof.

$$\langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{s}_t \rangle \stackrel{(1)}{\geq} \delta \max_{\mathbf{s} \in \mathcal{D}} \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{s} \rangle \quad (82)$$

$$= \delta \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t \rangle + \delta \max_{\mathbf{s} \in \mathcal{D}} \langle -\nabla f(\mathbf{x}_t), \mathbf{s} \rangle \quad (83)$$

$$= \delta (\langle \nabla f(\mathbf{x}_t), \mathbf{x}_t \rangle + \sigma_{\mathcal{D}}(-\nabla f(\mathbf{x}_t))) \quad (84)$$

$$= \delta (f(\mathbf{x}_t) + \underbrace{f^*(\nabla f(\mathbf{x}_t)) + \sigma_{\mathcal{D}}(-\nabla f(\mathbf{x}_t))}_{=-\psi(\nabla f(\mathbf{x}_t))}) = \delta (f(\mathbf{x}_t) - \psi(\nabla f(\mathbf{x}_t))) \quad (85)$$

where the first identity uses the definition of \mathbf{s}_t , the second one the definition of convex conjugate and the last one is a consequence of the Fenchel-Young identity. We recall $\sigma_{\mathcal{D}}$ is the support function of \mathcal{D} . \square

Theorem 2. Let f be convex, \mathbf{x}_t denote the iterate generated by any of the proposed FW variants (*AdaFW*, *AdaAFW*, *AdaPFW*) after t iterations, with $N_t \geq 1$, and let \mathbf{u}_t be defined recursively as $\mathbf{u}_0 = \nabla f(\mathbf{x}_0)$, $\mathbf{u}_{t+1} = (1 - \xi_t)\mathbf{u}_t + \xi_t \nabla f(\mathbf{x}_t)$, where $\xi_t = 2/(\delta N_t + 2)$ if t is a good step and $\xi_t = 0$ otherwise. Then we have:

$$h_t \leq f(\mathbf{x}_t) - \psi(\mathbf{u}_t) \leq \frac{2\bar{L}_t \text{diam}(\mathcal{A})^2}{\delta^2 N_t + \delta} + \frac{2(1 - \delta)}{\delta^2 N_t^2 + \delta N_t} (f(\mathbf{x}_0) - \psi(\mathbf{u}_0)) = \mathcal{O}\left(\frac{1}{\delta^2 t}\right). \quad (86)$$

Proof. The proof is structured as follows. First, we derive a bound for the case that k is a good step. Second, we derive a bound for the case that k is a bad step. Finally, we add over all iterates to derive the claimed bound.

Case 1: k is a good step:

By Lemma 7, we have the following sequence of inequalities, valid for all $\xi_t \in [0, 1]$:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \xi_k g_k + \frac{\xi_k^2 L_k}{2} \|\mathbf{d}_k\|^2 \quad (87)$$

$$\leq f(\mathbf{x}_k) - \xi_k \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{s}_k \rangle + \frac{\xi_k^2 L_k}{2} \|\mathbf{d}_k\|^2 \quad (88)$$

$$= (1 - \delta \xi_k) f(\mathbf{x}_k) + \delta \xi_k \psi(\nabla f(\mathbf{x}_k)) + \frac{\xi_k^2 L_k}{2} \|\mathbf{d}_k\|^2, \quad (89)$$

where the second inequality follows from the definition of g_k (this is an equality for *AdaFP* but an inequality for the other variants) and the last identity from Lemma 8.

We now introduce the auxiliary variable σ_k . This is defined recursively as $\sigma_0 = \psi(\nabla f(\mathbf{x}_k))$, $\sigma_{k+1} = (1 - \delta \xi_k) \sigma_k + \delta \xi_k \psi(\nabla f(\mathbf{x}_k))$. Subtracting σ_{k+1} from both sides of the previous inequality gives

$$f(\mathbf{x}_{k+1}) - \sigma_{k+1} \leq (1 - \delta \xi_k) [f(\mathbf{x}_k) - \sigma_k] + \frac{\xi_k^2 L_k}{2} \|\mathbf{s}_k - \mathbf{x}_k\|^2 \quad (90)$$

Let $\xi_k = 2/(\delta N_k + 2)$ and $a_k \stackrel{\text{def}}{=} \frac{1}{2}((N_k - 2)\delta + 2)((N_k - 1)\delta + 2)$. With these definitions, we have the following trivial identities that we will use soon:

$$a_{k+1}(1 - \delta \xi_k) = \frac{1}{2}((N_k - 2)\delta + 2)((N_k - 1)\delta + 2) = a_k \quad (91)$$

$$a_{k+1} \frac{\xi_k^2}{2} = \frac{((N_k - 1)\delta + 2)}{(N_k \delta + 2)} \leq 1 \quad (92)$$

where in the first inequality we have used that k is a good step and so $N_{k+1} = N_k + 1$.

Multiplying (90) by a_{k+1} we have

$$a_{k+1} (f(\mathbf{x}_{k+1}) - \sigma_{k+1}) \leq a_{k+1} (1 - \delta \xi_k) [f(\mathbf{x}_k) - \sigma_k] + \frac{L_k}{2} \|\mathbf{s}_k - \mathbf{x}_k\|^2 \quad (93)$$

$$\stackrel{(91)}{=} a_k [f(\mathbf{x}_k) - \sigma_k] + \frac{L_k}{2} \|\mathbf{s}_k - \mathbf{x}_k\|^2 \quad (94)$$

$$\leq a_k [f(\mathbf{x}_k) - \sigma_k] + L_k \text{diam}(\mathcal{A})^2 \quad (95)$$

Case 2: k is a bad step:

Lemma 2 with $\xi_k = 0$ guarantees that the objective function is non-increasing, i.e., $f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k)$. By construction of σ_k we have $\sigma_{k+1} = \sigma_k$, and so adding both multiplied by a_{k+1} we obtain

$$a_{k+1} (f(\mathbf{x}_{k+1}) - \sigma_{k+1}) \leq a_{k+1} (f(\mathbf{x}_k) - \sigma_k) \quad (96)$$

$$= a_k (f(\mathbf{x}_k) - \sigma_k), \quad (97)$$

where in the last identity we have used that its a bad step and so $a_{k+1} = a_k$.

Final: combining cases and adding over iterates:

We can combine (95) and (97) into the following inequality:

$$a_{k+1} (f(\mathbf{x}_k) - \sigma_k) - a_k (f(\mathbf{x}_k) - \sigma_k) \leq L_k \text{diam}(\mathcal{A})^2 \mathbb{1}\{k \text{ is a good step}\}, \quad (98)$$

where $\mathbb{1}\{\text{condition}\}$ is 1 if condition is verified and 0 otherwise.

Adding this inequality from 0 to $t - 1$ gives

$$a_t(f(\mathbf{x}_t) - \sigma_t) \leq \sum_{k=0}^{t-1} L_k Q_A^2 \mathbb{1}\{k \text{ is a good step}\} + a_0(f(\mathbf{x}_0) - \sigma_0) \quad (99)$$

$$= N_t \bar{L}_t \text{diam}(\mathcal{A})^2 + (1 - \delta)(2 - \delta)(f(\mathbf{x}_0) - \sigma_0) \quad (100)$$

Finally, dividing both sides by a_t (note that $a_t > 0$ for $N_t \geq 1$) and using $(2 - \delta) \leq 2$ we obtain

$$f(\mathbf{x}_t) - \sigma_t \leq \frac{2N_t}{((N_t - 2)\delta + 2)((N_t - 1)\delta + 2)} \bar{L}_t Q_A^2 \quad (101)$$

$$+ \frac{4(1 - \delta)}{((N_t - 2)\delta + 2)((N_t - 1)\delta + 2)} (f(\mathbf{x}_0) - \sigma_0) \quad (102)$$

We will now use the inequalities $(N_t - 2)\delta + 2 \geq N_t\delta$ and $(N_t - 1)\delta + 2 \geq N_t\delta + 1$ for the terms in the denominator to obtain

$$f(\mathbf{x}_t) - \sigma_t \leq \frac{2\bar{L}_t Q_A^2}{\delta^2 N_t + \delta} + \frac{4(1 - \delta)}{\delta_t^2 N_t^2 + \delta N_t} (f(\mathbf{x}_0) - f(\mathbf{x}^*)) . \quad (103)$$

In order to prove the claimed bound we just need to prove the bound $-\psi(\mathbf{u}_t) \leq -\sigma_t$. We will prove this by induction. For $t = 0$ we have $\psi(\mathbf{u}_t) = \sigma_t$ by definition and so the bound is trivially verified. Suppose its true for t , then for $t + 1$ we have

$$-\psi(\mathbf{u}_{t+1}) = -\psi((1 - \xi_t)\mathbf{u}_t + \xi_t \nabla f(\mathbf{x}_t)) \quad (104)$$

$$\leq -(1 - \xi_t)\psi(\mathbf{u}_t) - \xi_t \psi(\nabla f(\mathbf{x}_t)) \quad (105)$$

$$\leq -(1 - \xi_t)\sigma_t - \xi_t \psi(\nabla f(\mathbf{x}_t)) \quad (106)$$

$$= -\sigma_{t+1} \quad (107)$$

where the first inequality is true by convexity of $-\psi$ and the second one by the induction hypothesis. Using this bound in (103) yields the desired bound

$$f(\mathbf{x}_t) - \psi(\mathbf{u}_t) \leq \frac{2\bar{L}_t Q_A^2}{\delta^2 N_t + \delta} + \frac{4(1 - \delta)}{\delta_t^2 N_t^2 + \delta N_t} [f(\mathbf{x}_0) - \psi(\nabla f(\mathbf{x}_0))] \quad (108)$$

We will now prove the bound $h_t \leq f(\mathbf{x}_t) - \psi(\mathbf{u}_t)$. Let \mathbf{u}^* be an arbitrary maximizer of ψ . Then by duality we have that $f(\mathbf{x}^*) = \psi(\mathbf{u}^*)$ and so

$$f(\mathbf{x}_t) - \psi(\mathbf{u}_t) = f(\mathbf{x}_t) - f(\mathbf{x}^*) + \psi(\mathbf{u}^*) - \psi(\mathbf{u}_t) \geq f(\mathbf{x}_t) - f(\mathbf{x}^*) = h_t \quad (109)$$

Finally, the $\mathcal{O}(\frac{1}{\delta_t})$ rate comes from bounding the number of good steps from (9), for which we have $1/N_t \leq \mathcal{O}(1/t)$, and bounding the Lipschitz estimate by a contant (Proposition 2). \square

Appendix D.2 Matching Pursuit

Lemma 9. Let \mathbf{s}_t be as defined in *AdaMP*, $R_{\mathcal{B}}$ be the level set radius defined as

$$R_{\mathcal{B}} = \max_{\substack{\mathbf{x} \in \text{lin}(\mathcal{A}) \\ f(\mathbf{x}) \leq f(\mathbf{x}_0)}} \|\mathbf{x} - \mathbf{x}^*\|_{\mathcal{B}} , \quad (110)$$

and \mathbf{x}^* be any solution to (OPT-MP). Then we have

$$\langle -\nabla f(\mathbf{x}_t), \mathbf{s}_t \rangle \geq \frac{\delta}{\max\{\mathcal{R}_{\mathcal{B}}, 1\}} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \quad (111)$$

Proof. By definition of atomic norm we have

$$\frac{\mathbf{x}_t - \mathbf{x}^*}{\|\mathbf{x}_t - \mathbf{x}^*\|_{\mathcal{B}}} \in \text{conv}(\mathcal{B}) \quad (112)$$

Since $f(\mathbf{x}_t) \leq f(\mathbf{x}_0)$, which is a consequence of sufficient decrease condition (Eq. (72)), we have that $R_{\mathcal{B}} \geq \|\mathbf{x}_t - \mathbf{x}^*\|_{\mathcal{B}}$ and so $\zeta \stackrel{\text{def}}{=} \|\mathbf{x}_t - \mathbf{x}^*\|_{\mathcal{B}}/R_{\mathcal{B}} \leq 1$. By symmetry of \mathcal{B} we have that

$$\frac{\mathbf{x}_t - \mathbf{x}^*}{R_{\mathcal{B}}} = \zeta \frac{\mathbf{x}_t - \mathbf{x}^*}{\|\mathbf{x}_t - \mathbf{x}^*\|_{\mathcal{B}}} + (1 - \zeta)\mathbf{0} \in \text{conv}(\mathcal{B}) . \quad (113)$$

We will now use this fact to bound the original expression. By definition of \mathbf{s}_t we have

$$\langle -\nabla f(\mathbf{x}_t), \mathbf{s}_t \rangle \stackrel{(5)}{\geq} \delta \max_{\mathbf{s} \in \mathcal{B}} \langle -\nabla f(\mathbf{x}_t), \mathbf{s} \rangle \quad (114)$$

$$\stackrel{(113)}{\geq} \frac{\delta}{R_{\mathcal{B}}} \langle -\nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x}^* \rangle \quad (115)$$

$$\geq \frac{\delta}{R_{\mathcal{B}}} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \quad (116)$$

where the last inequality follows by convexity. □

Theorem 3. Let f be convex, \mathbf{x}^* be an arbitrary solution to (OPT-MP) and let $R_{\mathcal{B}}$ the level set radius:

$$R_{\mathcal{B}} = \max_{\substack{\mathbf{x} \in \text{lin}(\mathcal{A}) \\ f(\mathbf{x}) \leq f(\mathbf{x}_0)}} \|\mathbf{x} - \mathbf{x}^*\|_{\mathcal{B}} . \quad (117)$$

If we denote by \mathbf{x}_t the iterate generated by AdaMP after $t \geq 1$ iterations and $\beta = \delta/R_{\mathcal{B}}$, then we have:

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{2\bar{L}_t \text{radius}(\mathcal{A})^2}{\beta^2 t + \beta} + \frac{2(1 - \beta)}{\beta^2 t^2 + \beta t} h_0 = \mathcal{O}\left(\frac{1}{\beta^2 t}\right) . \quad (118)$$

Proof. Let \mathbf{x}^* be an arbitrary solution to (OPT-MP). Then by Lemma 2, we have the following sequence of inequalities, valid for all $\xi_t \geq 0$:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \xi_k \langle -\nabla f(\mathbf{x}_k), \mathbf{s}_k \rangle + \frac{\xi_k^2 L_k}{2} \|\mathbf{s}_k\|^2 \quad (119)$$

$$\leq f(\mathbf{x}_k) - \xi_k \frac{\delta}{R_{\mathcal{B}}} [f(\mathbf{x}_k) - f(\mathbf{x}^*)] + \frac{\xi_k^2 L_k}{2} \|\mathbf{s}_k\|^2 , \quad (120)$$

where the second inequality follows from Lemma 9.

Subtracting $f(\mathbf{x}^*)$ from both sides of the previous inequality gives

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq \left(1 - \frac{\delta}{R_{\mathcal{B}}} \xi_k\right) [f(\mathbf{x}_k) - f(\mathbf{x}^*)] + \frac{\xi_k^2 L_k}{2} \|\mathbf{s}_k\|^2 . \quad (121)$$

Let $\beta = \delta/R_{\mathcal{B}}$ and $\xi_k = 2/(\beta k + 2)$ and $a_k \stackrel{\text{def}}{=} \frac{1}{2}((k-2)\beta + 2)((k-1)\beta + 2)$. With these definitions, we have the following trivial results:

$$a_{k+1}(1 - \beta \xi_k) = \frac{1}{2}((k-2)\beta + 2)((k-1)\beta + 2) = a_k \quad (122)$$

$$a_{k+1} \frac{\xi_k^2}{2} = \frac{((k-1)\beta + 2)}{(k\beta + 2)} \leq 1 . \quad (123)$$

Multiplying (121) by a_{k+1} we have

$$a_{k+1}(f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*)) \leq a_{k+1}(1 - \beta\xi_k)[f(\mathbf{x}_k) - f(\mathbf{x}^*)] + \frac{L_k}{2}\|\mathbf{s}_k\|^2 \quad (124)$$

$$\stackrel{(91)}{=} a_k[f(\mathbf{x}_k) - f(\mathbf{x}^*)] + \frac{L_k}{2}\|\mathbf{s}_k\|^2 \quad (125)$$

$$\leq a_k[f(\mathbf{x}_k) - f(\mathbf{x}^*)] + L_t \text{radius}(\mathcal{A})^2 \quad (126)$$

Adding this last inequality from 0 to $t-1$ gives

$$a_t(f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \sum_{k=0}^{t-1} L_k \text{radius}(\mathcal{A})^2 + a_0(f(\mathbf{x}_0) - \beta_0) \quad (127)$$

$$= t\bar{L}_t \text{diam}(\mathcal{A})^2 + (1 - \delta)(2 - \delta)(f(\mathbf{x}_0) - \beta_0) \quad (128)$$

Finally, dividing both sides by a_t (note that $a_1 = 2 - \beta \geq 1$ and so a_t is strictly positive for $t \geq 1$), and using $(2 - \delta) \leq 2$ we obtain

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{2t}{((t-2)\beta + 2)((t-1)\beta + 2)} \bar{L}_t \text{radius}(\mathcal{A})^2 \quad (129)$$

$$+ \frac{4(1 - \beta)}{((t-2)\beta + 2)((t-1)\beta + 2)} (f(\mathbf{x}_0) - \beta_0) \quad (130)$$

We will now use the inequalities $(t-2)\beta + 2 \geq t\beta$ and $(t-1)\beta + 2 \geq t\beta + 1$ to simplify the terms in the denominator. With this we obtain to obtain

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{2\bar{L}_t \text{radius}(\mathcal{A})^2}{\beta^2 N_t + \beta} + \frac{4(1 - \beta)}{\beta_t^2 N_t^2 + \beta N_t} (f(\mathbf{x}_0) - f(\mathbf{x}^*)) , \quad (131)$$

which is the desired bound. \square

Appendix E Proofs of convergence for strongly convex objectives

The following proofs depend on some definitions of geometric constants, which are defined in [Appendix A](#) as well as two crucial lemmas from ([Lacoste-Julien and Jaggi, 2015](#), Appendix C).

Appendix E.1 Frank-Wolfe variants

We are now ready to present the convergence rate of the adaptive Frank–Wolfe variants. As we did in [Appendix C](#), although the original proof combines the rates for FW variants and MP, the proof will be split into two, in which we prove separately the linear convergence rates for [AdaAFW](#) and [AdaPFW](#) (Theorem 4.A) and [AdaMP](#) (Theorem 4.B).

Theorem 4.A. *Let f be μ -strongly convex. Then for each good step we have the following geometric decrease:*

$$h_{t+1} \leq (1 - \rho_t)h_t, \quad (132)$$

with

$$\rho_t = \frac{\mu\delta^2}{4L_t} \left(\frac{\text{PWidth}(\mathcal{A})}{\text{diam}(\mathcal{D})} \right)^2 \quad \text{for AdaAFW} \quad (133)$$

$$\rho_t = \min \left\{ \frac{\delta}{2}, \delta^2 \frac{\mu}{L_t} \left(\frac{\text{PWidth}(\mathcal{A})}{\text{diam}(\mathcal{D})} \right)^2 \right\} \quad \text{for AdaPFW} \quad (134)$$

Note. In the main paper we provided the simplified bound $\rho_t = \frac{\mu}{4L_t} \left(\frac{\text{PWidth}(\mathcal{A})}{\text{diam}(\mathcal{A})} \right)^2$ for both algorithms [AdaAFW](#) and [AdaPFW](#) for simplicity. It is easy to see that the bound for [AdaPFW](#) above can be trivially bounded by this quantity by noting that $\delta^2 \leq \delta$ and that μ/L_t and $\text{PWidth}(\mathcal{A})/\text{diam}(\mathcal{D})$ are necessarily smaller than 1.

Proof. The structure of this proof is similar to that of ([Lacoste-Julien and Jaggi, 2015](#), Theorem 8). We begin by upper bounding the suboptimality h_t . Then we derive a lower bound on $h_{t+1} - h_t$. Combining both we arrive at the desired geometric decrease.

Upper bounding h_t

Assume \mathbf{x}_t is not optimal, ie $h_t > 0$. Then we have $\langle -\nabla f(\mathbf{x}_t), \mathbf{x}^* - \mathbf{x}_t \rangle > 0$. Using the definition of the geometric strong convexity bound and letting $\bar{\gamma} \stackrel{\text{def}}{=} \gamma(\mathbf{x}_t, \mathbf{x}^*)$ we have

$$\frac{\bar{\gamma}^2}{2} \mu_f^A \leq f(\mathbf{x}^*) - f(\mathbf{x}_t) + \langle -\nabla f(\mathbf{x}_t), \mathbf{x}^* - \mathbf{x}_t \rangle \quad (135)$$

$$= -h_t + \bar{\gamma} \langle -\nabla f(\mathbf{x}_t), \mathbf{s}_f(\mathbf{x}_t) - \mathbf{v}_f(\mathbf{x}_t) \rangle \quad (136)$$

$$\leq -h_t + \bar{\gamma} \langle -\nabla f(\mathbf{x}_t), \mathbf{s}_t - \mathbf{v}_t \rangle \quad (137)$$

$$= -h_t + \bar{\gamma} q_t, \quad (138)$$

where $q_t \stackrel{\text{def}}{=} \langle -\nabla f(\mathbf{x}_t), \mathbf{s}_t - \mathbf{v}_t \rangle$. For the last inequality we have used the definition of $\mathbf{v}_f(\mathbf{x})$ which implies $\langle f(\mathbf{x}_t), \mathbf{v}_f(\mathbf{x}_t) \rangle \leq \langle \nabla f(\mathbf{x}_t), \mathbf{v}_t \rangle$ and the fact that $\mathbf{s}_t = \mathbf{s}_f(\mathbf{x}_t)$. Therefore

$$h_t \leq -\frac{\bar{\gamma}^2}{2} \mu_f^A + \bar{\gamma} q_t, \quad (139)$$

which can always be upper bounded by taking $\bar{\gamma} = \mu^{-1} q_t$ (since this value of $\bar{\gamma}$ maximizes the expression on the right hand side of the previous inequality) to arrive at

$$h_t \leq \frac{q_t^2}{2\mu_f^A} \quad (140)$$

$$\leq \frac{q_t^2}{2\mu\Delta^2}, \quad (141)$$

with $\Delta \stackrel{\text{def}}{=} \text{PWidth}(\mathcal{A})$ and where the last inequality follows from Lemma 1.

Lower bounding progress $h_t - h_{t+1}$.

Let G be defined as $G = 1/2$ for AdaAFW and $G = 1$ for AdaPFW. We will now prove that for both algorithms we have

$$\langle -\nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle \geq \delta G q_t. \quad (142)$$

For AdaAFW, by the way the direction \mathbf{d}_t is chosen on Line 6, we have the following sequence of inequalities:

$$\begin{aligned} 2\langle -\nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle &\geq \langle -\nabla f(\mathbf{x}_t), \mathbf{d}_t^{FW} \rangle + \langle -\nabla f(\mathbf{x}_t), \mathbf{d}_t^A \rangle \\ &\geq \delta \langle -\nabla f(\mathbf{x}_t), \mathbf{s}_t - \mathbf{x}_t \rangle + \delta \langle -\nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{v}_t \rangle \\ &= \delta \langle -\nabla f(\mathbf{x}_t), \mathbf{s}_t - \mathbf{v}_t \rangle \\ &= \delta q_t, \end{aligned}$$

For AdaPFW, since $\mathbf{d}_t = \mathbf{s}_t - \mathbf{v}_t$, it follows from the definition of q_t that $\langle -\nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle \geq \delta q_t$.

We split the rest of the analysis into three cases: $\gamma_t < \gamma_t^{\max}$, $\gamma_t = \gamma_t^{\max} \geq 1$ and $\gamma_t = \gamma_t^{\max} < 1$. We prove a geometric descent in the first two cases. In the case where $\gamma_t = \gamma_t^{\max} < 1$ (a bad step) we show that the number of bad steps is bounded.

Case 1: $\gamma_t < \gamma_t^{\max}$:

By Lemma 2, we have

$$f(\mathbf{x}_{t+1}) = f(\mathbf{x}_t + \gamma_t \mathbf{d}_t) \leq f(\mathbf{x}_t) + \min_{\eta \in [0, \gamma_t^{\max}]} \left\{ \eta \langle \nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle + \frac{L_t \eta^2}{2} \|\mathbf{d}_t\|^2 \right\} \quad (143)$$

Because $\gamma_t < \gamma_t^{\max}$ and since the expression inside the minimization term (143) is a convex function of η , the minimizer is unique and it coincides with the minimum of the unconstrained problem. Hence we have

$$\min_{\eta \in [0, \gamma_t^{\max}]} \left\{ \eta \langle \nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle + \frac{L_t \eta^2}{2} \|\mathbf{d}_t\|^2 \right\} = \min_{\eta \geq 0} \left\{ \eta \langle \nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle + \frac{L_t \eta^2}{2} \|\mathbf{d}_t\|^2 \right\} \quad (144)$$

Replacing in (2), our bound becomes

$$f(\mathbf{x}_{t+1}) = f(\mathbf{x}_t + \gamma_t \mathbf{d}_t) \leq f(\mathbf{x}_t) + \min_{\eta \geq 0} \left\{ \eta \langle \nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle + \frac{L_t \eta^2}{2} \|\mathbf{d}_t\|^2 \right\} \quad (145)$$

$$\leq f(\mathbf{x}_t) + \min_{\eta \geq 0} \left\{ \eta \langle \nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle + \frac{L_t \eta^2}{2} M^2 \right\} \quad (146)$$

$$\leq f(\mathbf{x}_t) + \eta \langle \nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle + \frac{L_t \eta^2}{2} M^2, \quad \forall \eta \geq 0 \quad (147)$$

where the second inequality comes from bounding $\|\mathbf{d}_t\|$ by $M \stackrel{\text{def}}{=} \text{diam}(\mathcal{D})$. Subtracting $f(\mathbf{x}^*)$ from both sides and rearranging we have

$$h_t - h_{t+1} \geq \eta \langle -\nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle - \frac{1}{2} \eta^2 L_t M^2, \quad \forall \eta \geq 0. \quad (148)$$

Using the gap inequality (142) our lower bound becomes

$$h_t - h_{t+1} \geq \eta \delta G q_t - \frac{1}{2} \eta^2 L_t M^2, \quad \forall \eta \geq 0. \quad (149)$$

Noting that the lower bound in (149) is a concave function of η , we maximize the bound by selecting $\eta^* = (L_t M^2)^{-1} \delta G q_t$. Plugging η^* into the bound in (149) and then using the strong convexity bound (141) we have

$$h_t - h_{t+1} \geq \frac{\mu G^2 \Delta^2 \delta^2}{L_t M^2} h_t \implies h_{t+1} \leq \left(1 - \frac{\mu G^2 \Delta^2 \delta^2}{L_t M^2} \right) h_t. \quad (150)$$

Then we have geometric convergence with rate $1 - \rho$ where $\rho = (4L_t M^2)^{-1} \mu \Delta^2 \delta^2$ for [AdaAFW](#) and $\rho = (L_t M^2)^{-1} \mu \Delta^2 \delta^2$ for [AdaPFW](#).

Case 2: $\gamma_t = \gamma_t^{\max} \geq 1$

By Lemma 2 and the gap inequality (142), we have

$$h_t - h_{t+1} = f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}) \geq \eta \delta G q_t - \frac{1}{2} \eta^2 L_t M^2, \quad \forall \eta \leq \gamma_t^{\max}. \quad (151)$$

Since the lower bound (151) is true for all $\eta \leq \gamma_t^{\max}$, we can maximize the bound with $\eta^* = \min\{(L_t M^2)^{-1} \delta G q_t, \gamma_t^{\max}\}$. In the case when $\eta^* = (L_t M^2)^{-1} \delta G q_t$ we get the same bound as we do in (150) and hence have linear convergence with rate $1 - \rho$ where $\rho = (4L_t M^2)^{-1} \mu \Delta^2 \delta^2$ for [AdaAFW](#) and $\rho = (L_t M^2)^{-1} \mu \Delta^2 \delta^2$ for [AdaPFW](#). If $\eta^* = \gamma_t^{\max}$ then this implies $L_t M^2 \leq \delta G q_t$. Since γ_t^{\max} is assumed to be greater than 1 and the bound holds for all $\eta \leq \gamma_t^{\max}$ we have in particular that it holds for $\eta = 1$ and hence

$$h_t - h_{t+1} \geq \delta G q_t - \frac{1}{2} L_t M^2 \quad (152)$$

$$\geq \delta G q_t - \frac{\delta G q_t}{2} \quad (153)$$

$$\geq \frac{\delta G h_t}{2}, \quad (154)$$

where in the second line we use the inequality $L_t M^2 \leq \delta G q_t$ and in the third we use the inequality $h_t \leq q_t$ which is an immediate consequence of convexity of f . Then we have

$$h_{t+1} \leq (1 - \rho) h_t, \quad (155)$$

where $\rho = \delta/4$ for [AdaAFW](#) and $\rho = \delta/2$ for [AdaPFW](#). Note by Proposition 1 and the fact $\mu \leq L_t$ we have $\delta/4 \geq (4L_t M^2)^{-1} \mu \Delta^2 \delta^2$.

Case 3: $\gamma_t = \gamma_t^{\max} < 1$ (bad step)

In this case, we have either a drop or swap step and can make no guarantee on the progress of the algorithm (drop and swap are defined in Appendix B). For [AdaAFW](#), $\gamma_t = \gamma_t^{\max} < 1$ is a drop step. From lines 6–9 of [AdaAFW](#) we can make the following distinction of cases. In case of a FW step, then $\mathcal{S}_{t+1} = \{\mathbf{s}_t\}$ and $\gamma_t = \gamma_t^{\max} = 1$, otherwise $\mathcal{S}_{t+1} = \mathcal{S}_t \cup \{\mathbf{s}_t\}$. In case of an Away step, $\mathcal{S}_{t+1} = \mathcal{S}_t \setminus \{\mathbf{v}_t\}$ if $\gamma_t = \gamma_t^{\max} < 1$, otherwise $\mathcal{S}_{t+1} = \mathcal{S}_t$. Note a drop step can only occur at an Away step. For [AdaPFW](#), $\gamma_t = \gamma_t^{\max} < 1$ will be a drop step when $\mathbf{s}_t \in \mathcal{S}_t$ and will be a swap step when $\mathbf{s}_t \notin \mathcal{S}_t$.

Even though at these bad steps we do not have the same geometric decrease, Lemma 5 yields that the sequence $\{h_t\}$ is a non-increasing sequence, i.e., $h_{t+1} \leq h_t$. Since we are guaranteed a geometric decrease on steps that are not bad steps, the bounds on the number of bad steps of Eq. (9) is sufficient to conclude that [AdaAFW](#) and [AdaPFW](#) exhibit a global linear convergence. \square

Appendix E.2 Matching Pursuit

We start by proving the following lemma, which will be crucial in the proof of the Adaptive MP's linear convergence rate.

Lemma 10. *Suppose that \mathcal{A} is a non-empty compact set and that f is μ -strongly convex. Let $\nabla_{\mathcal{B}} f(\mathbf{x})$ denote the orthogonal projection of $\nabla f(\mathbf{x})$ onto $\text{lin}(\mathcal{B})$. Then for all $\mathbf{x}^* - \mathbf{x} \in \text{lin}(\mathcal{A})$, we have*

$$f(\mathbf{x}^*) \geq f(\mathbf{x}) - \frac{1}{2\mu \text{mDW}(\mathcal{B})^2} \|\nabla_{\mathcal{B}} f(\mathbf{x})\|_{\mathcal{B}^*}^2. \quad (156)$$

Proof. From [Locatello et al. \(2018, Theorem 6\)](#), we have that if f is μ -strongly convex, then

$$\mu_{\mathcal{B}} \stackrel{\text{def}}{=} \inf_{\mathbf{x}, \mathbf{y} \in \text{lin}(\mathcal{B}), \mathbf{x} \neq \mathbf{y}} \frac{2}{\|\mathbf{y} - \mathbf{x}\|_{\mathcal{B}}^2} [f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle] \quad (157)$$

is positive and verifies $\mu_{\mathcal{B}} \geq \text{mDW}(\mathcal{B})^2 \mu$. Replacing $\mathbf{y} = \mathbf{x} + \gamma(\mathbf{x}^* - \mathbf{x})$ in the definition above we have

$$f(\mathbf{x} + \gamma(\mathbf{x}^* - \mathbf{x})) \geq f(\mathbf{x}) + \gamma \langle \nabla f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle + \gamma^2 \frac{\mu_{\mathcal{B}}}{2} \|\mathbf{x}^* - \mathbf{x}\|_{\mathcal{B}}^2. \quad (158)$$

We can fix $\gamma = 1$ on the left hand side and since the expression on the right hand side is true for all γ , we minimize over γ to find $\gamma^* = -\langle \nabla f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle / \mu_{\mathcal{B}} \|\mathbf{x}^* - \mathbf{x}\|_{\mathcal{B}}^2$. Thus the lower bound becomes

$$f(\mathbf{x}^*) \geq f(\mathbf{x}) - \frac{1}{2\mu_{\mathcal{B}}} \frac{\langle \nabla f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle}{\|\mathbf{x}^* - \mathbf{x}\|_{\mathcal{B}}^2} \quad (159)$$

$$\geq f(\mathbf{x}) - \frac{1}{2\mu \text{mDW}(\mathcal{B})^2} \frac{\langle \nabla f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle}{\|\mathbf{x}^* - \mathbf{x}\|_{\mathcal{B}}^2} \quad (160)$$

$$= f(\mathbf{x}) - \frac{1}{2\mu \text{mDW}(\mathcal{B})^2} \frac{\langle \nabla_{\mathcal{B}} f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle}{\|\mathbf{x}^* - \mathbf{x}\|_{\mathcal{B}}^2} \quad (161)$$

$$\geq f(\mathbf{x}) - \frac{1}{2\mu \text{mDW}(\mathcal{B})^2} \|\nabla_{\mathcal{B}} f(\mathbf{x})\|_{\mathcal{B}^*}^2, \quad (162)$$

where the last inequality follows by $|\langle \mathbf{y}, \mathbf{z} \rangle| \leq \|\mathbf{y}\|_{\mathcal{B}^*} \|\mathbf{z}\|_{\mathcal{B}}$ \square

Theorem 4.B. (Convergence rate Adaptive MP) Let f be μ -strongly convex and suppose \mathcal{B} is a non-empty compact set. Then *AdaMP* verifies the following geometric decrease for each $t \geq 0$:

$$h_{t+1} \leq (1 - \delta^2 \rho_t) h_t, \quad \text{with } \rho_t = \frac{\mu}{L_t} \left(\frac{\text{mDW}(\mathcal{B})}{\text{radius}(\mathcal{B})} \right)^2, \quad (163)$$

where $\text{mDW}(\mathcal{B})$ the minimal directional width of \mathcal{B} .

Proof. By Lemma 2 and bounding $\|\mathbf{d}_t\|$ by $R = \text{radius}(\mathcal{B})$ we have

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) + \min_{\eta \in \mathbb{R}} \left\{ \eta \langle \nabla f(\mathbf{x}_t), \mathbf{s}_t \rangle + \frac{\eta^2 L_t R^2}{2} \right\} \quad (164)$$

$$= f(\mathbf{x}_t) - \frac{\langle \nabla f(\mathbf{x}_t), \mathbf{s}_t \rangle^2}{2L_t R^2} \quad (165)$$

$$\leq f(\mathbf{x}_t) - \delta^2 \frac{\langle \nabla f(\mathbf{x}_t), \mathbf{s}_t^* \rangle^2}{2L_t R^2} \quad (166)$$

where \mathbf{s}_t^* is any element such that $\mathbf{s}_t^* \in \arg \min_{\mathbf{s} \in \mathcal{B}} \langle \nabla f(\mathbf{x}_t), \mathbf{s} \rangle$ and the inequality follows from the optimality of min and the fact that $\langle \nabla f(\mathbf{x}_t), \mathbf{s}_t^* \rangle \leq 0$. Let $\nabla_{\mathcal{B}} f(\mathbf{x}_t)$ denote as in Lemma 10 the orthogonal projection of $\nabla f(\mathbf{x}_t)$ onto $\text{lin}(\mathcal{B})$. Then the previous inequality simplifies to

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \delta^2 \frac{\langle \nabla_{\mathcal{B}} f(\mathbf{x}_t), \mathbf{s}_t^* \rangle^2}{2L_t R^2}. \quad (167)$$

By definition of dual norm, we also have $\langle -\nabla_{\mathcal{B}} f(\mathbf{x}_t), \mathbf{s}_t^* \rangle = \|\nabla_{\mathcal{B}} f(\mathbf{x}_t)\|_{\mathcal{B}^*}^2$. Subtracting $f(\mathbf{x}^*)$ from both sides we obtain the upper-bound:

$$h_{t+1} \leq h_t - \delta^2 \frac{\|\nabla_{\mathcal{B}} f(\mathbf{x}_t)\|_{\mathcal{B}^*}^2}{2L_t R^2} \quad (168)$$

To derive the lower-bound, we use Lemma 10 with $\mathbf{x} = \mathbf{x}_t$ and see that

$$\|\nabla_{\mathcal{B}} f(\mathbf{x}_t)\|_{\mathcal{B}^*} \geq 2\mu \text{mDW}(\mathcal{B})^2 h_t \quad (169)$$

Combining the upper and lower bound together we have

$$h_{t+1} \leq \left(1 - \delta^2 \frac{\mu \text{mDW}(\mathcal{B})^2}{L_t R^2} \right) h_t, \quad (170)$$

which is the claimed bound. \square

Appendix F Experiments

In this appendix we give some details on the experiments which were omitted from the main text, as well as an extended set of results.

Appendix F.1 ℓ_1 -regularized logistic regression, Madelon dataset

For the first experiment, we consider an ℓ_1 -regularized logistic regression of the form

$$\arg \min_{\|\mathbf{x}\|_1 \leq \beta} \frac{1}{n} \sum_{i=1}^n \varphi(\mathbf{a}_i^T \mathbf{x}, \mathbf{b}_i) + \frac{\lambda}{2} \|\mathbf{x}\|_2^2, \quad (171)$$

where φ is the logistic loss. The linear subproblems in this case can be computed exactly ($\delta = 1$) and consists of finding the largest entry of the gradient. The regularization parameter λ is always set to $\lambda = \frac{1}{n}$.

We first consider the case in which the data $\mathbf{a}_i, \mathbf{b}_i$ is the Madelon dataset. Below are the curves objective suboptimality vs time for the different methods considered. The regularization parameter, denoted ℓ_1 ball radius in the figure, is chosen as to give 1%, 5% and 20% of non-zero coefficients (the middle figure is absent from the main text).

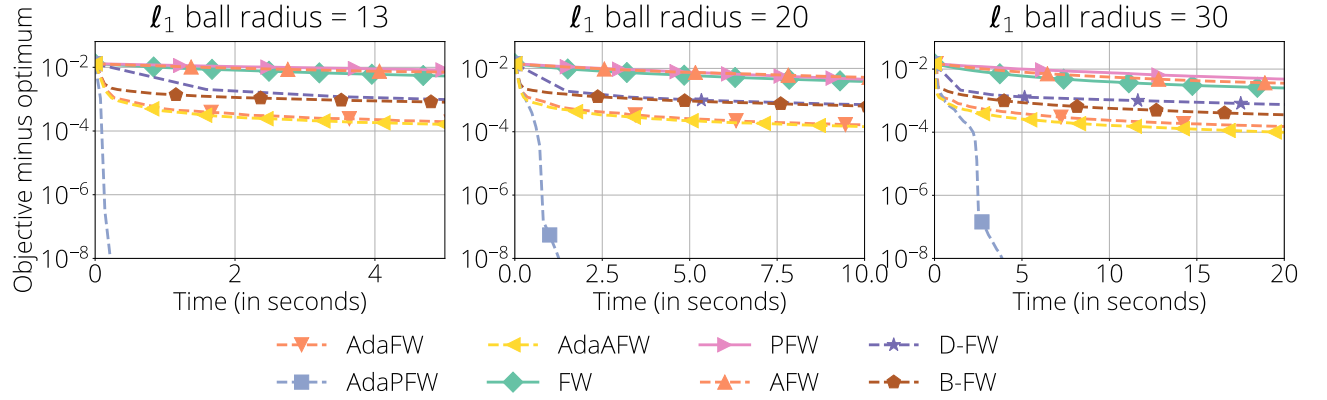


Figure 2: **Comparison of different FW variants.** Problem is ℓ_1 -regularized logistic regression and dataset is Madelon in the first, RCV1 in the second figure.

Appendix F.2 ℓ_1 -regularized logistic regression, RCV1 dataset

The second experiment is identical to the first one, except the madelon dataset is replaced by the larger RCV1 dataset. Below we display the results of the comparison in this dataset:

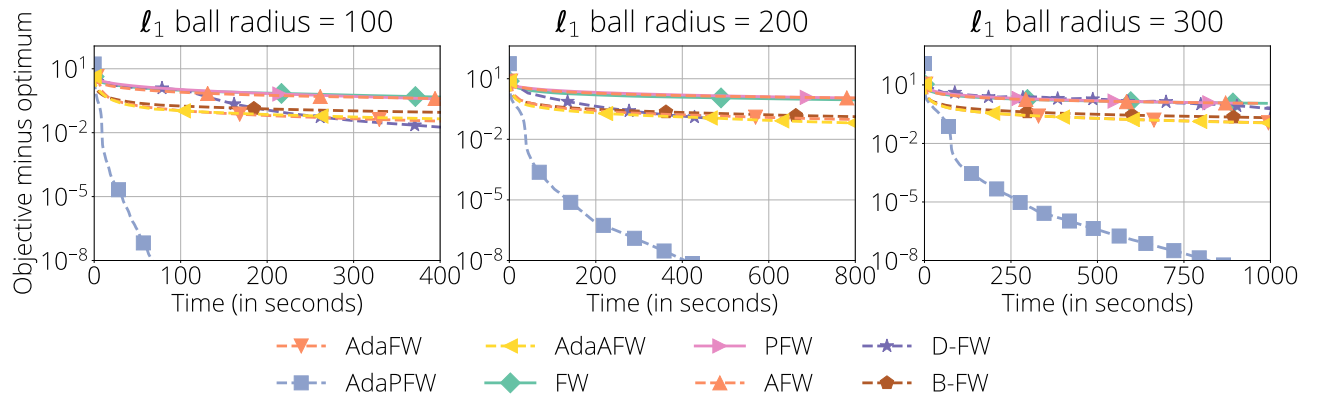


Figure 3: **Comparison of different FW variants.** Problem is ℓ_1 -regularized logistic regression and dataset is RCV1.

Appendix F.3 Nuclear norm-regularized Huber regression, MovieLens dataset

For the third experiment, we consider a collaborative filtering problem with the MovieLens 1M dataset (Harper and Konstan, 2015) as provided by the spotlight¹ Python package.

In this case the dataset consists of a sparse matrix \mathbf{A} representing the ratings for the different movies and users. We denote by \mathcal{I} the non-zero indices of this matrix. Then the optimization problem that we consider is the following

$$\arg \min_{\|\mathbf{X}\|_* \leq \beta} \frac{1}{n} \sum_{(i,j) \in \mathcal{I}} L_\xi(\mathbf{A}_{i,j} - \mathbf{X}_{i,j}), \quad (172)$$

where H_1 is the Huber loss, defined as

$$L_\xi(a) = \begin{cases} \frac{1}{2}a^2 & \text{for } |a| \leq \xi, \\ \xi(|a| - \frac{1}{2}\xi), & \text{otherwise.} \end{cases} \quad (173)$$

The Huber loss is a quadratic for $|a| \leq \xi$ and grows linearly for $|a| > \xi$. The parameter ξ controls this tradeoff and was set to 1 during the experiments.

We compared the variant of FW that do not require to store the active set on this problem (as these are the only competitive variants for this problem).

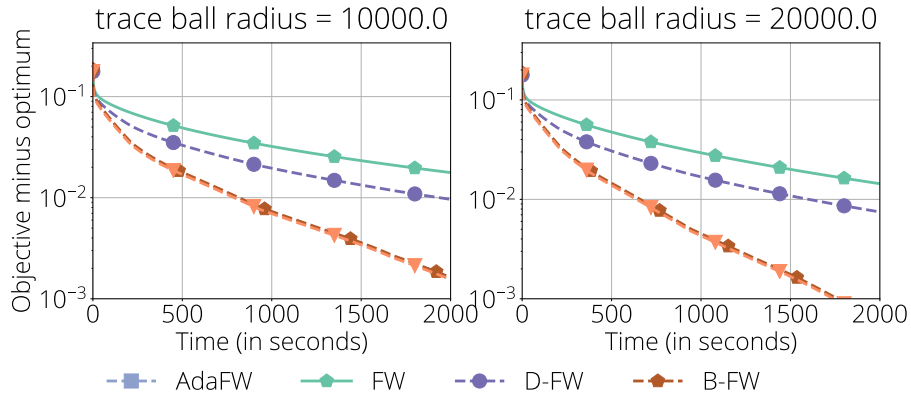


Figure 4: **Comparison of different FW variants.** Comparison of FW variants on the MovieLens 1M dataset.

¹<https://github.com/maciejkula/spotlight>