

1 **FIRST-ORDER METHODS FOR THE IMPATIENT:**
2 **SUPPORT IDENTIFICATION IN FINITE TIME**
3 **WITH CONVERGENT FRANK-WOLFE VARIANTS**

4 IMMANUEL M. BOMZE*, FRANCESCO RINALDI†, AND SAMUEL ROTA BULÒ‡

5 **Abstract.** In this paper, we focus on the problem of minimizing a non-convex function over the
6 unit simplex. We analyze two well-known and widely used variants of the Frank-Wolfe algorithm and
7 first prove global convergence of the iterates to stationary points both when using exact and Armijo
8 line search. Then we show that the algorithms identify the support in a finite number of iterations
9 (the identification result does not hold for the classic Frank-Wolfe algorithm). This, to the best of
10 our knowledge, is the first time a manifold identification property has been shown for such a class of
11 methods.

12 **Key words.** Surface Identification, Manifold Identification, Active Set, Finite Convergence

13 **AMS subject classifications.** 65K05, 90C06, 90C30

14 **1. Introduction.** The minimization of a (possibly non-convex) function over
15 the probability simplex is a problem arising in many different contexts like, e.g., ma-
16 chine learning, statistics and economics (see, e.g., [7, 11] for an overview of real-world
17 applications). When dealing with this kind of problems, Frank-Wolfe variants (see,
18 e.g., [16] and references therein) guarantee good scalability, thanks to the way they
19 handle the feasible set, and also give a sparse representation of the iterates, thus
20 offering a good alternative to projected gradient algorithms. Anyway, some may ar-
21 gue that projected gradient methods still represent the best choice in the considered
22 framework, since they can identify the sparsity pattern, i.e., the final set of non-zero
23 variables, in a finite number of iterations (under some specific assumptions). This fea-
24 ture is particularly useful if the solution of the problem is sparse and we just want to
25 find its support, since it means we do not need to run the algorithm until convergence.
26 It is also important when trying to speed-up a given algorithm. Indeed, after we iden-
27 tify the set of non-zero variables, we could simply apply some more sophisticated
28 Newton-like method over the lower-dimensional space those variables describe. Such
29 a feature may also help to develop suitable support identification/active-set strategies,
30 like the ones described in, e.g., [2, 4, 9, 10, 12, 14].

31 There exists a considerable number of papers analyzing support/active-set identifi-
32 cation properties of optimization methods. Bertsekas first showed in [1] that the
33 projected gradient method identifies the sparsity pattern in a finite number of iter-
34 ations when using non-negativity constraints. In [6] the authors showed that some
35 simple algorithms (including projected gradient) would, in a finite number of itera-
36 tions, identify the face of a polyhedral feasible region on which the solutions to an
37 optimization problem occur. These results were generalized in [24] to the case of
38 non-polyhedral convex sets. Analysis for nonconvex constraints is reported in [5, 15].

*ISOR, VCOR & ds:UniVie, Universität Wien, Austria (immanuel.bomze@univie.ac.at)

†Dipartimento di Matematica, Università di Padova, Italy (rinaldi@math.unipd.it)

‡Mapillary Research, Graz, Austria (samuel@mapillary.com)

39 The support identification property has also been established for other algorithms like
 40 certain coordinate descent and stochastic gradient methods [18, 25], proximal gradient
 41 methods (see, e.g., [19, 21]) and sequential minimal optimization methods for SVM
 42 training [22]. In [7], the problem of minimizing a convex function over the probability
 43 simplex is considered, and coreset-based results are reported for fully corrective ver-
 44 sions of some Frank-Wolfe variants. Recall that a coreset is a face of the simplex with
 45 the property that the minimum of the function on the face is a good approximate
 46 solution of the full problem. It is further important to remark that fully corrective
 47 algorithms heavily rely on the fact that a minimum of the function over a given face
 48 can be calculated at each iteration. Hence, those algorithms cannot be considered
 49 when dealing with non-convex problems.

50 In the present paper, we consider two well-known variants of the Frank-Wolfe
 51 algorithm, namely away-step Frank-Wolfe [23] and pairwise Frank-Wolfe [16, 20], and
 52 prove global convergence of their iterates to stationary points when using exact or
 53 Armijo line search (in the sense of characterizing all accumulation points of iterates
 54 by stationarity), and moreover global convergence for the full iteration sequence for
 55 the away-step variant. These results then enable us to prove support identification in a
 56 finite number of iterations for those algorithms. More specifically, when considering a
 57 convergent sequence (x^k) generated by one of those Frank-Wolfe variants, we have that
 58 it converges to a stationary point \bar{x} . Furthermore, we can be sure the iterates x^k will
 59 match the sparsity pattern of \bar{x} when k is sufficiently large (if strict complementarity
 60 holds at \bar{x}). This, to the best of our knowledge, is the first time that a support
 61 identification result is proved for Frank-Wolfe like algorithms.

62 This result is quite surprising if we take into account the fact that the classic
 63 Frank-Wolfe algorithm does not guarantee support identification in finite time. We
 64 will give some examples later on (see Section 4) where all iterates generated by the
 65 algorithm have full support (i.e., number of nonzero coordinates equal to the number
 66 of variables in the problem), and the limit point of the iterate sequence does not.

67 The paper is organized as follows. After a preliminary analysis of the problem in
 68 Section 2, we describe in depth the algorithmic framework in Section 3. In Section 4
 69 we establish global convergence and support identification property of the methods.
 70 Finally, in Section 5, we draw some conclusions.

71 **2. Preliminary Analysis of the Problem.** Denoting by $e = (1, \dots, 1)^\top$ the
 72 n -dimensional vector with all entries equal to one, the problem we consider here is
 73 the following:

$$74 \quad (2.1) \quad \min_{x \in \Delta} f(x)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $\Delta = \{x \in \mathbb{R}^n : e^\top x = 1, x \geq 0\}$ is the probability simplex. A
 class of C^2 -objective functions f including all quadratic functions will be considered
 in this paper. For any fixed $x \in \Delta$ and any feasible direction d (we will construct d
 such that always $[0, 1] \subseteq I_{\text{feas}}(x, d) := \{\alpha \in \mathbb{R} : x + \alpha d \in \Delta\}$ holds), denote by

$$\varphi_d^x(\alpha) = f(x + \alpha d), \quad \alpha \in I_{\text{feas}}(x, d)$$

75 with derivatives $\dot{\varphi}_d^x(\alpha) = d^\top \nabla f(x + \alpha d)$ and $\ddot{\varphi}_d^x(\alpha) = d^\top \nabla^2 f(x + \alpha d)d$.

76 We now give a key assumption on curvature of φ_d^x that will be needed to prove
 77 convergence of iterates (see Subsection 4.2). As we will see later on, this will guarantee
 78 that iteration is *homotopical* for the algorithms we analyze in the paper.

79 ASSUMPTION 2.1. Any φ_d^x is either concave or strictly convex over $I_{\text{feas}}(x, d)$.
 80 Furthermore, curvature should be bounded away from zero in the strictly convex case
 81 along descent directions: to be more precise, for all $x \in \Delta$ and all d with non-concave
 82 φ_d^x we ask existence of $\eta_d > 0$ such that

$$83 \quad (2.2) \quad \eta_d \leq \ddot{\varphi}_d^x(\alpha) \quad \text{for all } \alpha \in I_{\text{feas}}(x, d), \quad \text{if } \dot{\varphi}_d^x(0) < 0.$$

All quadratic functions $f(x) = x^\top Qx + c^\top x$, where Q is a possibly indefinite symmetric matrix, satisfy (2.2) with $\eta_d = d^\top Qd$ for all $x \in \Delta$. But many more functions f may meet Assumption 2.1, for example¹ the function $f(x) = c^\top x + \sqrt{x^\top Qx}$ for indefinite but strictly (co)positive Q (similar functions are used in volatility modeling). Then

$$[(x + \alpha d)^\top Q(x + \alpha d)]^{3/2} \ddot{\varphi}_d^x(\alpha) = (x^\top Qx)(d^\top Qd) - (d^\top Qx)^2$$

84 does not depend on α but can change sign with varying d .

85 For proving global convergence of the methods and support identification results
 86 (see Section 4), we need an essential global estimate following by continuity of $\nabla^2 f$
 87 over Δ (a set of diameter $\sqrt{2}$) made explicit in the following observation:

88 OBSERVATION 2.1. For all directions d with $\|d\| \leq \sqrt{2}$ and all $\alpha \in I_{\text{feas}}(x, d)$ we
 89 have bounded curvatures $\ddot{\varphi}_d^x(\alpha)$, or slightly more general $\|\nabla^2 f(x)\|_{\text{spec}} \leq K$ for all
 90 $x \in \Delta$ with the spectral norm $\|\cdot\|_{\text{spec}}$, implying

$$91 \quad (2.3) \quad |\dot{\varphi}_d^x(\alpha)| = |d^\top \nabla^2 f(x + \alpha d)d| \leq 2K \quad \text{for all } x \in \Delta, \quad \text{if } \|d\| \leq \sqrt{2}.$$

92 We further notice that minimizing a function $h(x)$ over a polytope P can be
 93 written as Problem (2.1). Let $V = [v_1, \dots, v_m] \in \mathbb{R}^{n \times m}$ be the matrix whose
 94 columns are the vertices of P . Since any point $y \in P$ can be expressed as a convex
 95 combination of the columns of V , the problem $\min\{h(y) : y \in P\}$ can be rewritten
 96 as the problem $\min\{f(x) = h(Vx) : x \in \Delta\}$. We note that

- 97 1. \bar{x} is a stationary point for f over Δ , cf. (3.1) below, if and only if $\bar{y} = V\bar{x}$ is
 98 a stationary point for h over P , i.e., satisfies the KKT conditions;
- 99 2. d is a descent direction for f at $x \in \Delta$ if Vd is one for h at $y = Vx \in P$; and
- 100 3. condition (2.2) carries over from h to f too, as $\nabla^2 f(x) = V^\top \nabla^2 h(Vx)V$.

101 **3. Frank-Wolfe Variants for Minimization over the Simplex.** In this sec-
 102 tion, we describe two well-known Frank-Wolfe variants that can be used to minimize
 103 a function over the probability simplex. In order to do that, we report below the
 104 generic scheme related to those iterative algorithms (see Algorithm 3.1). Beforehand
 105 we recall that $x^* \in \Delta$ is a stationary point for the problem (2.1) if and only if

$$106 \quad (3.1) \quad \nabla_r f(x^*) \geq \nabla f(x^*)^\top x^* \quad \text{for all } r \text{ with equality if } x_r^* > 0.$$

107 By construction, either the algorithm stops after finitely many iterations at a
 108 stationary point, or else the generated sequence takes infinitely many values in Δ as
 109 $f(x^{k+1}) < f(x^k)$.

¹We are grateful to Werner Schachinger who pointed to this in personal communication.

Algorithm 3.1 Line Search Algorithmic Scheme

- 1 Choose a point $x^0 \in \Delta$
 - 2 For $k = 0, 1, \dots$
 - 3 If x^k is a stationary point (3.1), then STOP
 - 4 Compute a feasible descent direction d^k at x^k
 - 5 Compute a stepsize $\alpha_k \in (0, 1]$ via line search for improving the objective
 - 6 Set $x^{k+1} = x^k + \alpha_k d^k$
 - 7 End for
-

110 **3.1. Frank-Wolfe type directions.** At every iteration k of Algorithm 3.1, we
 111 compute, at Step 4, a feasible descent search direction d^k that is used to generate the
 112 new iterate x^{k+1} . We describe here the different type of directions that can be used in
 113 Algorithm 3.1. We indicate the set of all indices related to the coordinates of vector
 114 x by $I = \{1, \dots, n\}$, and by $S_k = \{i \in I : x_i^k > 0\}$ we denote the support of x^k .

115 The Frank-Wolfe and the away-step directions (see, e.g., [13, 16]), computed in
 116 x^k are respectively:

$$117 \quad (3.2) \quad d_{FW}^{x^k} = e_{\hat{i}} - x^k, \quad \hat{i} \in \underset{i \in I}{\operatorname{Argmin}} \{ \nabla_i f(x^k) \};$$

118

$$119 \quad (3.3) \quad d_A^{x^k} = x^k - e_{\hat{j}}, \quad \hat{j} \in \underset{j \in S_k}{\operatorname{Argmax}} \{ \nabla_j f(x^k) \}.$$

120 We further indicate with $x_{\hat{j}}^k$ the \hat{j} -th coordinate of x^k , where \hat{j} is defined as in (3.3).
 121 Taking into account (3.2) and (3.3), we consider the following two search directions:

122 (AFW) Away-step Frank-Wolfe direction:

$$123 \quad d_{AFW}^{x^k} = \begin{cases} d_{FW}^{x^k}, & \text{if } \nabla f(x^k)^\top d_{FW}^{x^k} \leq \nabla f(x^k)^\top d_A^{x^k}, \\ \frac{x_{\hat{j}}^k}{1 - x_{\hat{j}}^k} d_A^{x^k}, & \text{otherwise.} \end{cases}$$

124 (PFW) Pairwise Frank-Wolfe direction:

$$125 \quad d_{PFW}^{x^k} = x_{\hat{j}}^k (d_{FW}^{x^k} + d_A^{x^k}) = x_{\hat{j}}^k (e_{\hat{i}} - e_{\hat{j}}),$$

126 where \hat{i} and \hat{j} are defined as in (3.2) and (3.3), respectively.

127 It is easy to verify that all above directions are strict descent directions, i.e., satisfy
 128 $\dot{\varphi}_d^x(0) = \nabla f(x)^\top d < 0$.

129 **3.2. Computation of the stepsize.** In the framework of Algorithm 3.1, given
 130 $x \in \Delta$ and a descent direction d at x , we aim at the largest (global) minimizer $\alpha_d^x > 0$
 131 of $\varphi_d^x(\alpha)$ over $(0, 1]$, i.e.

$$132 \quad (3.4) \quad \alpha_d^x := \max_{\alpha \in (0, 1]} \operatorname{Argmin} \varphi_d^x(\alpha).$$

Obviously, any global interior minimizer of φ in $(0, 1)$ satisfies the first-order condition

$$0 = \dot{\varphi}_d^x(\alpha_d^x) = \dot{\varphi}_d^x(0) + \alpha_d^x \ddot{\varphi}_d^x(\tilde{\alpha})$$

133 for some $\tilde{\alpha} \in [0, 1]$ depending on d and x . Hence, if $\ddot{\varphi}_d^x(\tilde{\alpha}) > 0$ we have

$$134 \quad (3.5) \quad \alpha_d^x = \frac{-\dot{\varphi}_d^x(0)}{\ddot{\varphi}_d^x(\tilde{\alpha})}.$$

135 **3.2.1. Exact and Armijo's line search.** Exact line search chooses, at a given
136 iteration k , the largest minimizer of $\varphi_{d^k}^{x^k}(\alpha)$ over $(0, 1]$, that is

$$137 \quad (3.6) \quad \alpha_k := \alpha_{d^k}^{x^k} \quad \text{defined as in (3.4) for } x = x^k \text{ and } d = d^k.$$

138 Unless the function $\varphi_{d^k}^{x^k}$ has some special structure (e.g., convexity/concavity),
139 determining the step size in (3.6) might be in general an expensive task. More prac-
140 tical strategies perform an inexact line search to identify the stepsize giving sufficient
141 reductions in the objective function at a minimal cost. A classic example is repre-
142 sented by the Armijo method (see, e.g., [3] and references therein). This method
143 iteratively shrinks the step size in order to guarantee a sufficient reduction of the
144 objective function. It represents a good way to replace exact line search in cases
145 when it gets too costly. In practice, we fix parameters $\delta \in (0, 1)$ and $\gamma \in (0, \frac{1}{2})$, and
146 start with maximal feasible stepsize equal to one. We then try steps $\alpha = \delta^m$ with
147 $m \in \mathbb{N}_0 = \{0, 1, 2, \dots\}$ until the sufficient decrease inequality

$$148 \quad (3.7) \quad f(x^k + \alpha d^k) \leq f(x^k) + \gamma \alpha \nabla f(x^k)^\top d^k$$

is satisfied, i.e., choose

$$m(x^k, d^k) := \min \{m \in \mathbb{N}_0 : (3.7) \text{ is satisfied for } \alpha = \delta^m\} < \infty$$

149 and set

$$150 \quad (3.8) \quad \alpha_k = \delta^{m(x^k, d^k)} \quad \text{as well as} \quad x^{k+1} = x^k + \alpha_k d^k.$$

151 Observe that under Assumption 2.1 on the curvature of φ_d^x , all stepsize variants we
152 discuss here have in common that always a full feasible step is taken unless in the
153 case of strictly convex φ_d^x where $\ddot{\varphi}_d^x(\alpha) > 0$ for all $\alpha \in [0, 1]$. So $\alpha_k < 1$ is possible
154 only if $\ddot{\varphi}_{d^k}^{x^k}(0) > 0$, for any strict descent direction d^k at x^k .

155 **3.2.2. Theoretical properties of line searches.** Now we prove that function
156 f reduces when moving from x^k to x^{k+1} , and that the sequence of the directional
157 derivatives along the search direction converges to zero when using the Armijo rule.
158 We will further see that a similar result also holds for the exact line search.

159 **PROPOSITION 3.1.** *Let (x^k) be the sequence generated by Algorithm 3.1 using*
160 *Armijo line search defined in (3.8), with any strict descent direction d^k satisfying*
161 *$\|d^k\| \leq \sqrt{2}$. Then we have*

- 162 (a) *if $x^{k+1} \neq x^k$, then $f(x^{k+1}) < f(x^k)$;*
163 (b) *if $x^{k+1} \neq x^k$ for all $k \in \mathbb{N}$, then $\lim_{k \rightarrow \infty} \nabla f(x^k)^\top d^k = 0$.*

Proof. We first notice that in a finite number of steps the Armijo line search finds a step satisfying condition (3.7). Then, due to the fact that d^k is such that $\nabla f(x^k)^\top d^k < 0$, we get that

$$f(x^{k+1}) < f(x^k).$$

164 Using again (3.7), we have

$$165 \quad (3.9) \quad f(x^k) - f(x^{k+1}) \geq \gamma \alpha_k |\nabla f(x^k)^\top d^k|.$$

166 Since $f(x^k)$ is monotonically decreasing and bounded in k , we can write

$$167 \quad (3.10) \quad \lim_{k \rightarrow \infty} \alpha_k |\nabla f(x^k)^\top d^k| = 0.$$

168 Let us consider, by contradiction, that (b) does not hold. In this case, due to the fact
169 that $\{\nabla f(x^k)^\top d^k\}$ is bounded, there exists an infinite subsequence k_j such that

$$170 \quad (3.11) \quad \lim_{j \rightarrow \infty} \nabla f(x^{k_j})^\top d^{k_j} = -\xi < 0,$$

171 with $\xi > 0$. Considering the limit in (3.10), we need to have

$$172 \quad (3.12) \quad \lim_{j \rightarrow \infty} \alpha_{k_j} = 0.$$

173 Using compactness of the feasible set Δ , we know that it is possible to get subsequence
174 (for ease of notation we again call it k_j) such that

$$175 \quad (3.13) \quad \lim_{k_j \rightarrow \infty} x^{k_j} = \hat{x} \quad \text{and} \quad \lim_{k_j \rightarrow \infty} d^{k_j} = \hat{d}.$$

176 Using continuity of the gradient, we thus can write

$$177 \quad (3.14) \quad \lim_{j \rightarrow \infty} \nabla f(x^{k_j})^\top d^{k_j} = \nabla f(\hat{x})^\top \hat{d} = -\xi < 0.$$

Taking into account (3.12), we in particular have for k_j sufficiently large

$$\alpha_{k_j} < 1.$$

178 Therefore

$$179 \quad (3.15) \quad f(x^{k_j} + \frac{\alpha_{k_j}}{\delta} d^{k_j}) - f(x^{k_j}) > \frac{\gamma \alpha_{k_j}}{\delta} \nabla f(x^{k_j})^\top d^{k_j}.$$

180 Using the mean value theorem we can replace the left hand side and write

$$181 \quad (3.16) \quad \frac{\alpha_{k_j}}{\delta} \nabla f(y^{k_j})^\top d^{k_j} > \gamma \frac{\alpha_{k_j}}{\delta} \nabla f(x^{k_j})^\top d^{k_j},$$

with $y^{k_j} = x^{k_j} + \theta_{k_j} \frac{\alpha_{k_j}}{\delta} d^{k_j}$ and $\theta_{k_j} \in (0, 1)$. Now, dividing by $\frac{\alpha_{k_j}}{\delta} > 0$ and taking into account that $y^{k_j} \rightarrow \hat{x}$ due to (3.12), we have

$$\nabla f(\hat{x})^\top \hat{d} \geq \gamma \nabla f(\hat{x})^\top \hat{d},$$

which finally gives us

$$\xi \leq \gamma \xi,$$

182 thus contradicting $\gamma < 1$ and proving that (b) holds. \square

Proposition 3.1 still holds when considering a stepsize $\bar{\alpha}_k \in (0, 1]$ satisfying the following inequality:

$$f(x^k + \bar{\alpha}_k d^k) \leq f(x^k + \alpha_k d^k),$$

183 where α_k is the stepsize obtained using the Armijo's rule. Indeed, if the above in-
 184 equality is satisfied, then (3.9) holds as well as the rest of the proof (see also Remark
 185 5 in [10]). Hence, we can easily get the following result:

186 **COROLLARY 3.2.** *Let (x^k) be the sequence of points in the feasible set Δ generated*
 187 *by Algorithm 3.1 using exact line search defined in (3.6), with any feasible descent*
 188 *direction d^k . Then we have*

- 189 (a) if $x^{k+1} \neq x^k$ then $f(x^{k+1}) < f(x^k)$;
 190 (b) if $x^{k+1} \neq x^k$ for all $k \in \mathbb{N}$, then $\lim_{k \rightarrow \infty} \nabla f(x^k)^\top d^k = 0$.

191 Summarizing Proposition 3.1 and Corollary 3.2, we get under the stepsize choice
 192 of (3.6) or (3.8) that

193 (3.17)
$$\dot{\varphi}_k(0) = \nabla f(x^k)^\top d^k \rightarrow 0 \quad \text{as } k \rightarrow \infty.$$

194 **4. Convergence results.** To clarify terminology, let us stress we use a common
 195 one: "global convergence" means that we establish the stationarity property for *all*
 196 *accumulation points* of the sequence of iterates (x^k) , regardless whether or not there
 197 may be more than one accumulation point. By contrast, "iterates convergence" means
 198 convergence of the full sequence (x^k) . Under mild assumptions which generically are
 199 true² we can show that there is only one accumulation point (which then enjoys
 200 stationarity by the global convergence results), if the sequence is generated by the
 201 (AFW) rule.

202 **4.1. Global convergence analysis.** In this section, for every considered choice
 203 of the direction d^k , we establish global convergence to stationary points of the algorith-
 204 mic scheme described above. Since the arguments for the different stepsize choices
 205 vary slightly, we chose to split the treatment. However, the two search direction
 206 choices are treated simultaneously, in an attempt of being concise.

207 **THEOREM 4.1.** *Let (x^k) be a sequence generated by Algorithm 3.1 where*

- 208 • the search direction d^k is computed according to (AFW) or (PFW) rule;
 209 • the stepsize α_k is computed using the Armijo line search described in (3.8).

210 *Then, either an integer $\bar{k} \geq 0$ exists such that $x^{\bar{k}}$ is a stationary point for prob-*
 211 *lem (2.1), or else the sequence (x^k) is infinite and every limit point x^* of the sequence*
 212 *is a stationary point (3.1) for problem (2.1).*

Proof. We first consider the case when the algorithm stops after a finite number
 of iterations \bar{k} . This can only happen if condition at Step 3 of Algorithm 3.1 is
 satisfied, i.e., if no direction d_{AFW} can be chosen, which is the case if and only if $x^{\bar{k}}$
 is a stationary point.

Now we consider the case when the sequence (x^k) is infinite. Arguing by contradiction,
 assume that there is an i such that $\nabla_i f(x^*) < \nabla f(x^*)^\top x^*$. We again distinguish cases:
Case 1 [not needed for (PFW)]. There is a subsequence k_j along which $x^{k_j} \rightarrow x^*$ and

²namely that there are only finitely many stationary points of the problem (2.1).

$d^{k_j} = d_{FW}^{x^{k_j}} = e^{i_j} - x^{k_j}$ for all j , where e^i denotes the i th column of the $n \times n$ identity matrix (and $i_j \in \{1, \dots, n\}$ suitably chosen). Then

$$\begin{aligned} \dot{\varphi}_{k_j}(0) &= \nabla f(x^{k_j})^\top d_{FW}^{x^{k_j}} = \nabla_{i_j} f(x^{k_j}) - \nabla f(x^{k_j})^\top x^{k_j} \\ &\leq \nabla_{i_j} f(x^{k_j}) - \nabla f(x^{k_j})^\top x^{k_j} \rightarrow \nabla_{i_j} f(x^*) - \nabla f(x^*)^\top x^* < 0 \end{aligned}$$

213 as $j \rightarrow \infty$, contradicting (3.17).

214 *Case 2a.* There is a subsequence k_j along which $x^{k_j} \rightarrow x^*$ and such that there is an
215 $\eta > 0$ with

$$216 \quad (4.1) \quad x_{r_j}^{k_j} \geq \eta \text{ for all } j,$$

217 where in the (AFW) case

$$218 \quad (4.2) \quad d^{k_j} = \frac{x_{r_j}^{k_j}}{1 - x_{r_j}^{k_j}} d_A^{x^{k_j}} = \frac{x_{r_j}^{k_j}}{1 - x_{r_j}^{k_j}} (x^{k_j} - e^{r_j})$$

219 whereas in the (PFW) case

$$220 \quad (4.3) \quad d^{k_j} = x_{r_j}^{k_j} (d_{FW}^{x^{k_j}} + d_A^{x^{k_j}}) = x_{r_j}^{k_j} (e^{\tilde{r}_j} - e^{r_j})$$

with $e^{\tilde{r}_j}$ the Frank-Wolfe vertex and e^{r_j} the away-step vertex. Then in the (AFW) case, as $\frac{x_{r_j}^{k_j}}{1 - x_{r_j}^{k_j}} \geq \frac{\eta}{1 - \eta} > 0$ holds for all j , we arrive at

$$\begin{aligned} \frac{1-\eta}{\eta} \dot{\varphi}_{k_j}(0) &= \frac{1-\eta}{\eta} \nabla f(x^{k_j})^\top d_A^{x^{k_j}} \\ &\leq \nabla f(x^{k_j})^\top d_A^{x^{k_j}} \leq \nabla f(x^{k_j})^\top d_{FW}^{x^{k_j}} \leq \nabla_{i_j} f(x^{k_j}) - \nabla f(x^{k_j})^\top x^{k_j} \\ &\rightarrow \nabla_{i_j} f(x^*) - \nabla f(x^*)^\top x^* < 0 \text{ as } j \rightarrow \infty, \end{aligned}$$

again contradicting (3.17). Similarly in the (PFW) case, the contradiction is obtained via

$$\begin{aligned} \frac{1}{\eta} \dot{\varphi}_{k_j}(0) &= \frac{1}{\eta} \nabla f(x^{k_j})^\top (d_{FW}^{x^{k_j}} + d_A^{x^{k_j}}) \\ &\leq \nabla f(x^{k_j})^\top (d_{FW}^{x^{k_j}} + d_A^{x^{k_j}}) \leq \nabla f(x^{k_j})^\top d_{FW}^{x^{k_j}} \\ &\leq \nabla_{i_j} f(x^{k_j}) - \nabla f(x^{k_j})^\top x^{k_j} \rightarrow \nabla_{i_j} f(x^*) - \nabla f(x^*)^\top x^* < 0 \end{aligned}$$

221 as $j \rightarrow \infty$. Hence the only remaining possibility is now

222 *Case 2b.* **any** convergent subsequence $x^{k_j} \rightarrow x^*$ with limit x^* satisfies

$$223 \quad (4.4) \quad x_{r_j}^{k_j} \rightarrow 0 \text{ as } j \rightarrow \infty,$$

where eventually (4.2) or (4.3) holds. Irrespective of the chosen direction, at least one such sequence (s_j) exists by assumption that x^* is a limit point of (x^k) . Consider

this subsequence and their immediate successors $k_j = s_j + 1$. By (4.2) or (4.3), we know

$$\|x^{s_j+1} - x^{s_j}\| \leq \alpha_{s_j} \max \left\{ \left\| \frac{x_{r_j}^{k_j}}{1 - x_{r_j}^{k_j}} d_A^{x^{k_j}} \right\|, \|x_{r_j}^{k_j} (d_{FW}^{x^{k_j}} + d_A^{x^{k_j}})\| \right\} \leq \sqrt{2} \frac{x_{r_j}^{k_j}}{1 - x_{r_j}^{k_j}} \rightarrow 0$$

224 as $j \rightarrow \infty$, since $\|d_A^{x^{k_j}}\| = \|x^{k_j} - e^{r_j}\| \leq \text{diam}\Delta = \sqrt{2}$ and likewise $\|d_{FW}^{x^{k_j}} + d_A^{x^{k_j}}\| =$
 225 $\|e^{\bar{r}_j} - e^{r_j}\| \leq \text{diam}\Delta = \sqrt{2}$. Therefore also $x^{s_j+1} \rightarrow x^*$ as $j \rightarrow \infty$, and we may
 226 also consider (4.4) with (4.2) or (4.3) for the successor sequence $k_j = s_j + 1$. By
 227 suitable thinning (s_j) if necessary, we may and do assume that eventually $r_j = r$ for
 228 all j . Then $x_r^{s_j+1} = 0$ eventually holds because otherwise $\alpha_{s_j} < 1$ and Prop. A.2
 229 applies, contradicting (4.4). Applying our conclusion (4.4) with (4.2) or (4.3) now to
 230 $k_j = s_j + 1$, we see that also an away-step $x^{s_j+1} - e^h$ (or a PFW step involving e^h
 231 as away-step vertex) with $h \neq r$ is selected for $k = s_j + 1$ (if j is large enough) with
 232 the property (again, after suitable thinning) that $x_h^{s_j+1} \rightarrow 0$ as $j \rightarrow \infty$ but still we
 233 have, by construction of the away (or PFW) step, $x_r^{s_j+2} = 0$ for all large enough j .
 234 So again we have $x^{s_j+2} \rightarrow x^*$ as $j \rightarrow \infty$, hence an index $t \notin \{r, h\}$ would be chosen
 235 for the away step at $k = s_j + 2$, and repeating the argument less than n times, no
 236 choice for d_A would be left, which is absurd in view of the fact that the sequence is
 237 infinite, whence neither Case 1 nor Case 2a applies. So the theorem is proved. \square

238 We close this section by proving global convergence of the Algorithm 3.1 when
 239 using the exact line search for calculating the stepsize.

240 **THEOREM 4.2.** *Let (x^k) be a sequence generated by Algorithm 3.1 where*

- 241 \bullet *the search direction d^k is computed according to (AFW) or (PFW) rule;*
- 242 \bullet *the stepsize α_k is computed using the line search described in (3.6).*

243 *Then, either an integer $\bar{k} \geq 0$ exists such that $x^{\bar{k}}$ is a stationary point for prob-*
 244 *lem (2.1), or else the sequence (x^k) is infinite and every limit point x^* of the sequence*
 245 *is a stationary point (3.1) for problem (2.1).*

246 *Proof.* The proof is very similar to the one given for the Armijo line search. The
 247 only difference is in Case 2b where we yield a contradiction by applying Prop. A.1. \square

248 **4.2. Iterates convergence and support identification in finite time.** We
 249 start with a general observation, in particular applicable to (AFW) and (PFW) direc-
 250 tions. All we need is that the conclusions of Theorems 4.1 and 4.2 hold, namely that
 251 all accumulation points are stationary; with this property, any strict local minimizer
 252 which is isolated among all stationary points can be shown to attract all sequences
 253 generated by Algorithm 3.1 which start close enough to it. Conversely, if the limit
 254 point attracts all iterates starting close enough to it, then necessarily this must be an
 255 isolated stationary point and a strict local minimizer of f over Δ .

256 Note that the equivalence below holds irrespective whether or not there are non-
 257 strict local solutions to (2.1).

258 **THEOREM 4.3.** *Let Assumption 2.1 hold. Consider Algorithm 3.1 with any de-*
 259 *scend direction and any stepsize, such that all accumulation points of generated se-*
 260 *quences (x^k) are stationary. Then the following two statements on a stationary point*
 261 *$p \in \Delta$ are equivalent:*

- 262 (a) there is a p -neighbourhood $U \subseteq \Delta$ with no stationary point in $U \setminus \{p\}$, and
 263 $f(x) > f(p)$ for all $x \in U \setminus \{p\}$.
 264 (b) there is a p -neighbourhood $V \subseteq \Delta$ such that every sequence (x^k) starting at
 265 $x^0 \in V$ converges to p .

Proof. (a) \Rightarrow (b): Let $\varepsilon > 0$ be so small that $B := \{x \in \Delta : \|x - p\| \leq \varepsilon\} \subseteq U$ and define

$$\sigma := \min \{f(x) : x \in \Delta, \|x - p\| = \varepsilon\} - f(p) > 0.$$

266 Then $V = \{x \in \Delta : f(x) < f(p) + \sigma, \|x - p\| < \varepsilon\} \subset U$ is relatively open in Δ and
 267 contains p , so a neighbourhood of p in Δ . We claim that any sequence starting in V
 268 will remain there forever. Indeed, suppose $x^{k+1} \notin V$ but $x^k \in V$ for some k ; then by
 269 convexity or concavity of f along $\text{conv}(x^k, x^{k+1})$ we have

$$270 \quad (4.5) \quad f(\lambda x^{k+1} + (1 - \lambda)x^k) \leq f(x^k) < f(p) + \sigma \quad \text{for all } \lambda \in [0, 1],$$

271 so $x^{k+1} \notin V$ would imply $\|x^{k+1} - p\| \geq \varepsilon$ and hence $\|\lambda x^{k+1} + (1 - \lambda)x^k - p\| = \varepsilon$ for
 272 some $\lambda \in (0, 1]$, contradicting the definition of σ . By compactness, all accumulation
 273 points of (x^k) must lie in B and thus in U . Since all of them are stationary by as-
 274 sumption, there can only be one, namely p , which means that (b) holds.

275 (b) \Rightarrow (a): Choose $U := V$. By monotonicity and continuity, we have $f(p) =$
 276 $\inf_k f(x^k) < f(x^0)$ for all $x^0 \in U \setminus \{p\}$, a set which does not contain any station-
 277 ary point, as all sequences starting there have to converge to p by assumption.
 278 \square

279 We thus have shown that in our model, every strict local solution is isolated
 280 (among all alternative stationary points $\tilde{x} \in \Delta$), which generally is not the case.
 281 Inspection of the proof of Theorem 4.3 reveals that the only essential property is that
 282 the iteration is *homotopical*, i.e. that the inequality on the left-hand side in (4.5) holds.
 283 We can conclude that for all these homotopical iteration procedures, convergence to
 284 a saddle point is highly unlikely, which is in line with recent findings in this research
 285 area for other first-order methods, see, e.g. [17] and references therein. Note that most
 286 of these papers deal with smooth transition maps (which facilitate characterization of
 287 saddle points via the Jacobian matrix) while our transition maps lack even continuity.

288 Next we need an auxiliary observation which only applies to d_{AFW} :

289 LEMMA 4.4. *Let Assumption 2.1 hold. Let $\gamma = \lim_{k \rightarrow \infty} f(x^k) = \inf_{k \in \mathbb{N}} f(x^k)$ and as-*
 290 *sume $\gamma = f(e^i)$ for some $i \in I$. Consider a certain iteration counter k with $x^{k+1} \neq x^k$.*
 291 *Then the following implications hold for both stepsize choices (3.6) or (3.8):*

- 292 (a) if $d^k = d_{FW}^{x^k} = e^i - x^k$, then Algorithm 3.1 stops at $k + 1$;
 293 (b) if $x_i^k > 0$ and $d^k = \frac{x_i^k}{1 - x_i^k} d_A^{x^k} = \frac{x_i^k}{1 - x_i^k} (x^k - e^i)$, then $x_i^{k+1} = 0$.

294 *Proof.* (a) By construction and assumption, we have $0 \leq f(x^{k+1}) - \gamma \leq f(e^i) -$
 295 $\gamma = 0$, hence $x^{k+2} = x^{k+1}$ which is a stationary point, using Proposition 3.1(a) or
 296 Corollary 3.2(a).

297 (b) Suppose that $\alpha_k < 1$; then φ_k has to be strictly convex, and by smoothness,
 298 f has to be strictly convex over the whole interval $\text{conv}(x^{k+1}, e^i)$. But as assumed
 299 above, we have $f(e^i) = \gamma \leq f(x^{k+1}) < f(x^k)$ in contradiction to the fact that
 300 $x^k \in \text{conv}(x^{k+1}, e^i)$. So necessarily $\alpha_k = 1$ and therefore $x_i^{k+1} = 0$. \square

301 We proceed to establish a convergence result for the full sequence of iterates under
302 mild assumptions for the away-step Frank-Wolfe variant:

303 **THEOREM 4.5.** *Let Assumption 2.1 hold. Consider a sequence (x^k) generated*
304 *by Algorithm 3.1 with stepsize choice (3.6) or (3.8), and d_{AFW} as descent direction.*
305 *Suppose that (x^k) has finitely many accumulation points. Then it must converge:*
306 *there is a $p \in \Delta$ such that $x^k \rightarrow p$ as $k \rightarrow \infty$.*

Proof. The statement obviously needs a proof only if the sequence (x^k) is infinite. So suppose there are finitely many (pairs of) accumulation points, but at least two. Choose pairwise disjoint neighbourhoods around all of them and wait until all x^k lie exactly in one of these neighbourhoods if $k \geq k_0$. Then, arguing by contradiction, if x^k would not converge, there is a subsequence k_j with $k_1 \geq k_0$ such that $x^{k_j} \rightarrow p$ and the immediate successors $x^{k_j+1} \rightarrow q \neq p$ as $j \rightarrow \infty$ which implies $\bar{\alpha} := \inf_j \alpha_{k_j} > 0$. Now, by thinning (k_j) if necessary, we may and do assume that $\alpha_{k_j} \rightarrow \alpha_\infty > 0$ as $j \rightarrow \infty$, and that there is an $i \in I$ with $d^{k_j} = e^i - x^{k_j}$ for all j , or else $d^{k_j} = \frac{x_i^{k_j}}{1-x_i^{k_j}}(x^{k_j} - e^i)$ with $x_i^{k_j} > 0$ for all j . Moreover, in this case we even get $x_i^{k_j} > c$ for all j and a suitable constant $c > 0$ because of

$$0 < \|q - p\| = \lim_j \|x^{k_j+1} - x^{k_j}\| \leq \sqrt{2}\alpha_\infty \lim_j \frac{x_i^{k_j}}{1-x_i^{k_j}} = \sqrt{2}\alpha_\infty \frac{p_i}{1-p_i}.$$

307 Next suppose that eventually stepsize is smaller than one, and we are in the strictly
308 convex case. Then, employing (3.17) and

$$309 \quad (4.6) \quad f(x^{k+1}) - f(x^k) = \varphi_k(\alpha_k) - \varphi_k(0) = \alpha_k \left[\dot{\varphi}_k(0) + \frac{\alpha_k}{2} \ddot{\varphi}_k(\hat{\alpha}_k) \right],$$

310 we obtain

$$311 \quad (4.7) \quad \ddot{\varphi}_{k_j}(\hat{\alpha}_{k_j}) \rightarrow 0.$$

Furthermore by continuity we have for any $\alpha \in I_{\text{feas}}(p, e^i - p)$,

$$\ddot{\varphi}_{k_j}(\alpha) \rightarrow (e^i - p)^\top \nabla^2 f((1 - \alpha)p + \alpha e^i)(e^i - p),$$

in the FW case, and for any $\alpha \in I_{\text{feas}}(p, \mu(p - e^i))$,

$$\ddot{\varphi}_{k_j}(\alpha) \rightarrow \mu^2(p - e^i)^\top \nabla^2 f((1 + \alpha\mu)p - \alpha\mu e^i)(p - e^i),$$

312 with $\mu = \frac{p_i}{1-p_i}$, in the away step case. On the other hand, we can employ (4.7) in
313 all cases. Now by (2.2) in Assumption 2.1 for $x = (1 - \alpha)p + \alpha e^i$, $\alpha \in [0, 1]$, and by
314 choosing $d = e^i - p$ in the FW case or $d = \mu(p - e^i)$ in the away step case, we conclude
315 that f must be linear along $\text{conv}(p, e^i)$ with slope $\lim_j \dot{\varphi}_{k_j}(0) = 0$, so constant, and
316 $f(e^i) = f(p) = \inf_k f(x^k)$ results.

317 Now in case of the FW direction, Lemma 4.4(a) would then yield the absurd conclusion
318 that Algorithm 3.1 stops even at iteration $k_1 + 1$.

319 In case of the away direction, we conclude by Lemma 4.4(b) that $x_i^{k_j+1} = 0$. But since
320 $x_i^{k_j+1} > 0$, we must have a FW step $d^k = e^i - x^k$ for some $k \in \{k_j + 1, \dots, k_{j+1} - 1\}$.
321 Now we again invoke Lemma 4.4(a) to arrive at the contradiction that Algorithm 3.1
322 stops at iteration $k + 1$, using $f(e^i) = f(p) = \inf_k f(x^k)$.

323 So we are left with the case that the stepsize equals eventually one. But then the
 324 argument is even simpler: in the FW case, we stop at e^i , and in the away case we
 325 directly get $x_i^{k_j+1} = 0$ and, as argued just before, stop again at e^i at some iteration
 326 counter $k \in \{k_j + 1, \dots, k_{j+1} - 1\}$ as well. \square

327 As a corollary to Theorems 4.1, 4.2 and 4.5, we thus obtain a generic convergence
 328 result for the iterates generated by Algorithm 3.1 for the away-step variant:

329 **COROLLARY 4.6.** *Suppose that (2.1) admits only finitely many stationary points.*
 330 *Then any sequence (x^k) generated by Algorithm 3.1 with stepsize choice (3.6) or (3.8),*
 331 *and d_{AFW} as descent direction, must converge: there is a $p \in \Delta$ such that $x^k \rightarrow p$ as*
 332 *$k \rightarrow \infty$.*

Now we introduce three sets that will be useful when carrying out the analysis related to support identification in finite time. More specifically, we call

$$S_+(x) := \{i \in I \mid \nabla_i f(x) > x^\top \nabla f(x)\},$$

$$S_-(x) := \{i \in I \mid \nabla_i f(x) < x^\top \nabla f(x)\},$$

and

$$S_0(x) := \{i \in I \mid \nabla_i f(x) = x^\top \nabla f(x)\}.$$

333 We hence report the announced results on support identification in finite time;
 334 note that strict complementarity (again generically true) of the stationary point \bar{x}
 335 exactly means $\bar{S} = S_0(\bar{x})$ in below theorem; recall that $S_-(\bar{x}) = \emptyset$ by (3.1).

336 **THEOREM 4.7.** *Consider a convergent sequence of iterates (x^k) , with supports*
 337 *$S_k = S(x^k)$, generated by Algorithm 3.1 for the following specifications:*

- 338 • *the search direction d^k is computed according to (AFW) or (PFW) rule;*
- 339 • *the stepsize α_k is computed using the line search described in (3.6) or (3.8).*

Denote by $\bar{x} := \lim_{k \rightarrow \infty} x^k$ as well as $\bar{S} := \{i \in I : \bar{x}_i > 0\}$, so that by stationarity (3.1) of \bar{x} we have $\bar{S} \subseteq S_0(\bar{x})$. Then there is a finite \bar{k} such that

$$\bar{S} \subseteq S_k \subseteq S_0(\bar{x}) \quad \text{for all } k \geq \bar{k}.$$

Proof. We can assume that $x^k = e^i$, with $i \in I$, cannot happen infinitely often. Indeed, otherwise by Lemma 4.4 the algorithm would stop after a finite number of iterations. So, we assume that $x^k \neq e^i$ for k sufficiently large. Now, by continuity of the gradient, we can find an iterate such that both the following inclusions hold:

$$S_+(\bar{x}) \subseteq S_+(x^k) \quad \text{and} \quad \bar{S} \subseteq S_k.$$

From stationarity of \bar{x} we can further write $\bar{S} \subseteq S_0(\bar{x}) = I \setminus S_+(\bar{x})$. Hence, we have

$$S_-(x^k) \subseteq I \setminus S_+(x^k) \subseteq I \setminus S_+(\bar{x}) = S_0(\bar{x})$$

340 implying

$$341 \quad (4.8) \quad S_-(x^k) \subseteq S_0(\bar{x}).$$

342 We claim now that once

$$343 \quad (4.9) \quad S_k \subseteq S_0(\bar{x})$$

holds for some k , then (4.9) is guaranteed for all the following iterations. Indeed, either $S_{k+1} = S_k \cup \{i\}$ and $i \in S_-(x^k) \subseteq S_0(\bar{x})$ or else the support is a subset of the current support, i.e., $S_{k+1} \subseteq S_k$. By contradiction to (4.9), let us assume that, when k sufficiently large, the set $S_k \setminus S_0(\bar{x})$ is never empty. Again, by continuity of the gradient, we can choose a sufficiently large k_0 to ensure existence of a positive value $\varrho > 0$ such that

$$|\nabla f(x^{k_j})^\top (e^i - x^{k_j})| < \varrho \quad \text{for all } i \in S_0(\bar{x}) \quad \text{whenever } k \geq k_0;$$

and

$$\nabla f(x^k)^\top (e^r - x^k) > \varrho \quad \text{for all } r \in S_k \setminus S_0(\bar{x}) = S_k \cap S_+(\bar{x}) \quad \text{whenever } k \geq k_0.$$

Hence, for both direction variants (AFW) and (PFW), we have that $e^{r(k)}$ is chosen in the algorithm as away-step vertex for some $r(k) \in S_k \setminus S_0(\bar{x})$, if $k \geq k_0$. Further, due to the finiteness of I , by considering a suitable subsequence k_j we can assume

$$r(k_j) = r \in S_{k_j} \setminus S_0(\bar{x}) = S_{k_j} \cap S_+(\bar{x}).$$

344 By stationarity of \bar{x} we know $r \notin \bar{S}$, so the r -th coordinate of \bar{x} satisfies $\bar{x}_r = 0$.
 345 Eventually, $x_r^{k_j+1} = 0$ holds exactly because otherwise $\alpha_{k_j} < 1$ and Prop. A.1 or
 346 Prop. A.2 applies, contradicting $x_r^{k_j} \rightarrow \bar{x}_r = 0$. Repeating the same argument for all
 347 other indices in $S_k \setminus S_0(\bar{x})$, the result is proved. \square

348 As we pointed out in the introduction, the classic Frank-Wolfe algorithm does not
 349 guarantee support identification in finite time. Below we report an example where all
 350 iterates x^k have full support (i.e., $|S^k| = n$) and the point \bar{x} does not (i.e., $|\bar{S}| < n$).

EXAMPLE 4.1 (Bad behaviour of the Frank-Wolfe algorithm). *We consider problem (2.1) having a quadratic objective function*

$$f(x) = \frac{1}{2} x^\top Q x,$$

with

$$Q = \begin{bmatrix} 6 & 0 & 6 \\ 0 & 3 & 3 \\ 6 & 3 & 10 \end{bmatrix}.$$

351 *It is easy to see that the solution in this case is the global minimizer $\bar{x} = (\frac{1}{3} \frac{2}{3} 0)^\top$.*
 352 *If we choose as starting point $x^0 = (0.1 \ 0.3 \ 0.6)^\top$, the Frank-Wolfe algorithm will not*
 353 *be able to get an iterate with the same support as \bar{x} in finite time, neither via exact*
 354 *nor with Armijo line search [8].*

355 Moreover, the behaviour of this version may be even deceptive as the support of
 356 the iterates is eventually constant also for this algorithm; indeed, either the iterates
 357 coincide with a vertex e^i infinitely often, so that monotonicity would imply finite
 358 convergence to e^i . But this is the benevolent case. In the opposite case, eventually no
 359 vertex is hit exactly during iterations, so that supports must (weakly) increase with
 360 t . By finiteness it follows that they remain eventually constant, but, as the example
 361 shows, S^k may overestimate the correct support \bar{S} .

362 **5. Conclusions.** In this paper, we studied methods for solving minimization
 363 problems over the probability simplex. More specifically, we analyzed two variants
 364 of the Frank-Wolfe algorithm, namely away-step and pairwise Frank-Wolfe. We first
 365 proved convergence of the iterates to stationary points both when using exact and
 366 Armijo line search, and even convergence for the full sequence of iterates for the away-
 367 step variant, under mild regularity assumptions. Then we showed that both discussed
 368 variants algorithms guarantee support identification in finite time, a property shared
 369 by projected gradient methods. As a future development, it may be worth while to
 370 analyze conditions which allow to get explicit bounds on the number of iterations
 371 required to identify the support correctly.

372 Appendix A. Auxiliary results.

373 **PROPOSITION A.1.** *Let $(x^{s_j}) \rightarrow x^*$ as $j \rightarrow \infty$ be a convergent subsequence gen-*
 374 *erated by Algorithm 3.1 according to (AFW) or (PFW) rule, where we abbreviate*
 375 *$d_{FW}^j = d_{FW}^{s_j}$ and $d_A^j = d_A^{s_j}$. We assume that for some fixed r , we have for all j*

- 376 • $d_{AFW}^{s_j} = \frac{x_r^{s_j}}{1-x_r^{s_j}} d_A^j = \frac{x_r^{s_j}}{1-x_r^{s_j}} (x^{s_j} - e^r)$ or $d_{PFW}^{s_j} = x_r^{s_j} (e^{i_j} - e^r)$;
- 377 • the stepsize is computed using the line search described in (3.6) and satisfies
- 378 $\alpha_{s_j} < 1$;
- 379 • one of the following cases holds:
 - 380 1. there exists i such that $\nabla_i f(x^*) < \nabla f(x^*)^\top x^*$, or
 - 381 2. there exists $\varrho > 0$ such that $\nabla f(x^{s_j})^\top (e^r - x^{s_j}) > \varrho$.

382 Then $x_r^* > 0$.

383 *Proof.* Since $\alpha_{s_j} < 1$ we have that (3.5) holds for some $\tilde{\alpha}_{k_j} \in [0, 1]$. So we arrive
 384 via (2.3) at

$$\begin{aligned}
 385 \quad \frac{x_r^{s_j}}{1-x_r^{s_j}} &\geq \alpha_{s_j} \frac{x_r^{s_j}}{1-x_r^{s_j}} = \frac{-\dot{\varphi}_{s_j}(0)}{\dot{\varphi}_{s_j}(\tilde{\alpha}_{s_j})} \frac{x_r^{s_j}}{1-x_r^{s_j}} = \frac{-\nabla f(x^{s_j})^\top d_A^j}{[d_A^j]^\top \nabla^2 f(x^{s_j} + \tilde{\alpha}_{k_j} d^{s_j}) d_A^j} \\
 386 \quad &\geq \frac{-\nabla f(x^{s_j})^\top d_A^j}{2K} \geq \frac{-\nabla f(x^{s_j})^\top d_{FW}^j}{2K} \geq \frac{\nabla f(x^{s_j})^\top x^{s_j} - \nabla_i f(x^{s_j})}{2K} \\
 387 \quad &\rightarrow \frac{\nabla f(x^*)^\top x^* - \nabla_i f(x^*)}{2K} > 0 \text{ as } j \rightarrow \infty,
 \end{aligned}$$

388 if we assume that Case 1 holds, and we get the same inequality for $j \rightarrow \infty$ also in
 389 Case 2 since

$$390 \quad -\nabla f(x^{s_j})^\top d_A^j > \varrho > 0.$$

391 This implies $x_r^* > 0$ for the (AFW) rule and likewise

$$\begin{aligned}
 392 \quad x_r^{s_j} &\geq \alpha_{s_j} x_r^{s_j} = \frac{-\dot{\varphi}_{s_j}(0)}{\dot{\varphi}_{s_j}(\tilde{\alpha}_{s_j})} x_r^{s_j} = \frac{-\nabla f(x^{s_j})^\top [d_{FW}^j + d_A^j]}{[d_{FW}^j + d_A^j]^\top \nabla^2 f(x^{s_j} + \tilde{\alpha}_{k_j} d^{s_j}) [d_{FW}^j + d_A^j]} \\
 393 \quad &\geq \frac{-\nabla f(x^{s_j})^\top [d_{FW}^j + d_A^j]}{2K} \geq \frac{-\nabla f(x^{s_j})^\top d_{FW}^j}{2K} \geq \frac{\nabla f(x^{s_j})^\top x^{s_j} - \nabla_i f(x^{s_j})}{2K} \\
 394 \quad &\rightarrow \frac{\nabla f(x^*)^\top x^* - \nabla_i f(x^*)}{2K} > 0 \text{ as } j \rightarrow \infty
 \end{aligned}$$

395 proves the result in Case 1 with the (PFW) rule, and the same inequality for $j \rightarrow \infty$
 396 holds in Case 2 since

$$397 \quad -\nabla f(x^{s_j})^\top [d_{FW}^j + d_A^j] > -\nabla f(x^{s_j})^\top d_A^j > \varrho > 0. \quad \square$$

398 **PROPOSITION A.2.** *Let $(x^{s_j}) \rightarrow x^*$ as $j \rightarrow \infty$ be a convergent subsequence gen-*
 399 *erated by Algorithm 3.1 according to (AFW) or (PFW) rule, where we abbreviate*
 400 *$d_{FW}^j = d_{FW}^{x^{s_j}}$ and $d_A^j = d_A^{x^{s_j}}$. We assume that for some fixed r , we have for all j*

- 401 • $d_{AFW}^{x^{s_j}} = \frac{x_r^{s_j}}{1-x_r^{s_j}} d_A^j = \frac{x_r^{s_j}}{1-x_r^{s_j}} (x^{s_j} - e^r)$ or $d_{PFW}^{x^{s_j}} = x_r^{s_j} (e^{i_j} - e^r)$;
- 402 • the stepsize is computed using the Armijo line search described in (3.8) and
- 403 satisfies $\alpha_{s_j} < 1$;
- 404 • one of the following cases holds:
 - 405 1. there exists i such that $\nabla_i f(x^*) < \nabla f(x^*)^\top x^*$, or
 - 406 2. there exists $\varrho > 0$ such that $\nabla f(x^{s_j})^\top (e^r - x^{s_j}) > \varrho$.

407 Then $x_r^* > 0$.

Proof. We first notice that for any $\alpha \in [0, 1]$ and $k = s_j$, by (2.3) we can write

$$f(x^k + \alpha d^k) \leq f(x^k) + \alpha \nabla f(x^k)^\top d^k + \frac{\alpha^2 K}{2} \|d^k\|^2.$$

So the sufficient decrease condition (3.7) would be satisfied if

$$f(x^k) + \alpha \nabla f(x^k)^\top d^k + \frac{\alpha^2 K}{2} \|d^k\|^2 \leq f(x^k) + \gamma \alpha \nabla f(x^k)^\top d^k,$$

and the latter holds true if

$$\alpha \leq \alpha_k^{\max} := \frac{2(1-\gamma) |\nabla f(x^k)^\top d^k|}{K \|d^k\|^2}.$$

408 This gives us an interval $[0, \alpha_k^{\max}]$ of step sizes satisfying sufficient decrease. Now, if
 409 $\alpha_k < 1$ is chosen as Armijo step size, then either $\alpha_k > \alpha_k^{\max}$ or else $\alpha_k \in [0, \alpha_k^{\max}]$
 410 but then $\frac{\alpha_k}{\delta} > \alpha_k^{\max}$ as the step size $\alpha = \frac{\alpha_k}{\delta}$ shall violate (3.7) by definition (3.8). In
 411 both cases, we get

$$412 \quad \alpha_k > \delta \alpha_k^{\max}.$$

413 Now we consider the two different search directions, abbreviating $d_{FW}^j = d_{FW}^{x^{s_j}}$ and

414 $d_A^j = d_A^{x^{s_j}}$. For the (AFW) rule, Case 1, we can write

$$\begin{aligned}
415 \quad & \frac{x_r^{s_j}}{1 - x_r^{s_j}} \geq \alpha_{s_j} \frac{x_r^{s_j}}{1 - x_r^{s_j}} > \delta \alpha_{s_j}^{\max} \frac{x_r^{s_j}}{1 - x_r^{s_j}} \\
416 \quad & = \frac{-\nabla f(x^{s_j})^\top d_A^j}{\|d_A^j\|^2} \frac{2\delta(1-\gamma)}{K} \geq \frac{-\nabla f(x^{s_j})^\top d_A^j}{2} \frac{2\delta(1-\gamma)}{K} \\
417 \quad & \geq \frac{-\nabla f(x^{s_j})^\top d_{FW}^j}{2} \frac{2\delta(1-\gamma)}{K} \\
418 \quad & \geq \frac{\delta(1-\gamma)}{K} [\nabla f(x^{s_j})^\top x^{s_j} - \nabla_i f(x^{s_j})] \\
419 \quad & \rightarrow \frac{\delta(1-\gamma)}{K} [\nabla f(x^*)^\top x^* - \nabla_i f(x^*)] > 0 \text{ as } j \rightarrow \infty,
\end{aligned}$$

420 and the same inequality for $j \rightarrow \infty$ can be obtained in Case 2 since

$$421 \quad -\nabla f(x^{s_j})^\top d_A^j > \varrho.$$

422 This implies $x_r^* > 0$. Similarly, for (PFW), Case 1, we can write

$$\begin{aligned}
423 \quad & x_r^{s_j} \geq \alpha_{s_j} x_r^{s_j} > \delta \alpha_{s_j}^{\max} x_r^{s_j} \\
424 \quad & = \frac{-\nabla f(x^{s_j})^\top [d_A^j + d_{FW}^j]}{\|d_A^j + d_{FW}^j\|^2} \frac{2\delta(1-\gamma)}{K} \geq -\nabla f(x^{s_j})^\top [d_A^j + d_{FW}^j] \frac{\delta(1-\gamma)}{K} \\
425 \quad & \geq -\nabla f(x^{s_j})^\top d_{FW}^{s_j} \frac{\delta(1-\gamma)}{K} \\
426 \quad & \geq \frac{\delta(1-\gamma)}{K} [\nabla f(x^{s_j})^\top x^{s_j} - \nabla_i f(x^{s_j})] \\
427 \quad & \rightarrow \frac{\delta(1-\gamma)}{K} [\nabla f(x^*)^\top x^* - \nabla_i f(x^*)] > 0 \text{ as } j \rightarrow \infty,
\end{aligned}$$

428 and the same inequality holds for $j \rightarrow \infty$ in Case 2 since

$$429 \quad -\nabla f(x^{s_j})^\top [d_{FW}^j + d_A^j] > -\nabla f(x^{s_j})^\top d_A^j > \varrho > 0. \quad \square$$

430

REFERENCES

- 431 [1] D. P. BERTSEKAS, *On the Goldstein-Levitin-Polyak gradient projection method*, IEEE Trans-
432 actions on automatic control, 21 (1976), pp. 174–184.
433 [2] D. P. BERTSEKAS, *Projected newton methods for optimization problems with simple constraints*,
434 SIAM J. Control Optim., 20 (1982), pp. 221–246.
435 [3] D. P. BERTSEKAS, *Nonlinear programming*, Athena scientific Belmont, 1999.
436 [4] E. G. BIRGIN AND J. M. MARTÍNEZ, *Large-scale active-set box-constrained optimization method*
437 *with spectral projected gradients*, Comput. Optim. Appl., 23 (2002), pp. 101–125.

- 438 [5] J. BURKE, *On the identification of active constraints ii: The nonconvex case*, SIAM Journal
439 on Numerical Analysis, 27 (1990), pp. 1081–1102.
- 440 [6] J. V. BURKE AND J. J. MORÉ, *On the identification of active constraints*, SIAM Journal on
441 Numerical Analysis, 25 (1988), pp. 1197–1211.
- 442 [7] K. L. CLARKSON, *Coresets, sparse greedy approximation, and the Frank-Wolfe algorithm*, ACM
443 Transactions on Algorithms (TALG), 6 (2010), p. 63.
- 444 [8] A. CRISTOFARI, M. DE SANTIS, S. LUCIDI, AND F. RINALDI, *New active-set frank-wolfe variants
445 for minimization over the simplex and the ℓ_1 -ball*, arXiv preprint arXiv:1703.07761, (2017).
- 446 [9] A. CRISTOFARI, M. DE SANTIS, S. LUCIDI, AND F. RINALDI, *A two-stage active-set algorithm
447 for bound-constrained optimization*, J. Optim. Theory Appl., 172 (2017), pp. 369–401.
- 448 [10] A. CRISTOFARI, M. DE SANTIS, S. LUCIDI, AND F. RINALDI, *An active-set algorithmic framework
449 for non-convex optimization problems over the simplex*, arXiv preprint arXiv:1703.07761v2,
450 (2018).
- 451 [11] E. DE KLERK, *The complexity of optimizing over a simplex, hypercube or sphere: a short
452 survey*, Central European Journal of Operations Research, 16 (2008), pp. 111–125.
- 453 [12] M. DE SANTIS, G. DI PILLO, AND S. LUCIDI, *An active set feasible method for large-scale mini-
454 mization problems with bound constraints*, Computational Optimization and Applications,
455 53 (2012), pp. 395–423.
- 456 [13] J. GUÉLAT AND P. MARCOTTE, *Some comments on Wolfe’s away step*, Mathematical Program-
457 ming, 35 (1986), pp. 110–119.
- 458 [14] W. W. HAGER AND H. ZHANG, *A new active set algorithm for box constrained optimization*,
459 SIAM J. Optim., 17 (2006), pp. 526–557.
- 460 [15] W. HARE AND A. S. LEWIS, *Identifying active constraints via partial smoothness and prox-
461 regularity*, Journal of Convex Analysis, 11 (2004), pp. 251–266.
- 462 [16] S. LACOSTE-JULIEN AND M. JAGGI, *On the global linear convergence of Frank-Wolfe optimiza-
463 tion variants*, in NIPS 2015 - Advances in Neural Information Processing Systems, 2015.
- 464 [17] J. D. LEE, I. PANAGEAS, G. PILIOURAS, M. SIMCHOWITZ, M. I. JORDAN, AND B. RECHT, *First-
465 order methods almost always avoid saddle points*, arXiv preprint arXiv:1710.07406, (2017).
- 466 [18] S. LEE AND S. J. WRIGHT, *Manifold identification in dual averaging for regularized stochastic
467 online learning*, Journal of Machine Learning Research, 13 (2012), pp. 1705–1744.
- 468 [19] R. MIFFLIN AND C. SAGASTIZÁBAL, *Proximal points are on the fast track*, Journal of Convex
469 Analysis, 9 (2002), pp. 563–580.
- 470 [20] B. MITCHELL, V. F. DEM’YANOV, AND V. MALOZEMOV, *Finding the point of a polyhedron
471 closest to the origin*, SIAM Journal on Control, 12 (1974), pp. 19–26.
- 472 [21] J. NUTINI, M. SCHMIDT, AND W. HARE, *“Active-set complexity” of proximal gradient: How
473 long does it take to find the sparsity pattern?*, Optimization Letters, (2018), [https://doi.
474 org/10.1007/s11590-018-1325-z](https://doi.org/10.1007/s11590-018-1325-z).
- 475 [22] J. SHE AND M. SCHMIDT, *Linear convergence and support vector identification of sequential
476 minimal optimization*, 10th NIPS Workshop on Optimization for Machine Learning, (2017).
- 477 [23] P. WOLFE, *Convergence theory in nonlinear programming*, Integer and nonlinear programming,
478 (1970), pp. 1–36.
- 479 [24] S. J. WRIGHT, *Identifiable surfaces in constrained optimization*, SIAM Journal on Control and
480 Optimization, 31 (1993), pp. 1063–1079.
- 481 [25] S. J. WRIGHT, *Accelerated block-coordinate relaxation for regularized optimization*, SIAM Jour-
482 nal on Optimization, 22 (2012), pp. 159–186.