
Frank-Wolfe Algorithms for Saddle Point Problems

Gauthier Gidel

INRIA - SIERRA Project-team
École normale supérieure, Paris

Tony Jebara

Department of Computer Science
Columbia University, NYC

Simon Lacoste-Julien

Department of CS & OR (DIRO)
Université de Montréal, Montréal

Abstract

We extend the Frank-Wolfe (FW) optimization algorithm to solve constrained smooth convex-concave saddle point (SP) problems. Remarkably, the method only requires access to linear minimization oracles. Leveraging recent advances in FW optimization, we provide the first proof of convergence of a FW-type saddle point solver over polytopes, thereby partially answering a 30 year-old conjecture. We also survey other convergence results and highlight gaps in the theoretical underpinnings of FW-style algorithms. Motivating applications without known efficient alternatives are explored through structured prediction with combinatorial penalties as well as games over matching polytopes involving an exponential number of constraints.

1 Introduction

The Frank-Wolfe (FW) optimization algorithm (Frank and Wolfe, 1956), also known as the conditional gradient method (Demjanov and Rubinov, 1970), is a first-order method for smooth constrained optimization over a compact set. It has recently enjoyed a surge in popularity thanks to its ability to cheaply exploit the structured constraint sets appearing in machine learning applications (Jaggi, 2013; Lacoste-Julien and Jaggi, 2015). A known forte of FW is that it only requires access to a *linear minimization oracle* (LMO) over the constraint set, i.e., the ability to minimize linear functions over the set, in contrast to projected gradient methods which require the minimization of *quadratic* functions or other nonlinear functions. In this paper, we extend the applicability of the FW algorithm to solve the following convex-concave saddle point problems:

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}(\mathbf{x}, \mathbf{y}), \quad (1)$$

with only access to $\text{LMO}(\mathbf{r}) \in \arg \min_{\mathbf{s} \in \mathcal{X} \times \mathcal{Y}} \langle \mathbf{s}, \mathbf{r} \rangle$,

where \mathcal{L} is a smooth (with L -Lipschitz continuous gradient) *convex-concave function*, that is, $\mathcal{L}(\cdot, \mathbf{y})$ is convex for all $\mathbf{y} \in \mathcal{Y}$ and $\mathcal{L}(\mathbf{x}, \cdot)$ is concave for all $\mathbf{x} \in \mathcal{X}$. We also assume that $\mathcal{X} \times \mathcal{Y}$ is a convex compact set such that its LMO is cheap to compute. A *saddle point solution* to (1) is a pair $(\mathbf{x}^*, \mathbf{y}^*) \in \mathcal{X} \times \mathcal{Y}$ (Hiriart-Urruty and Lemaréchal, 2013, VII.4) such that: $\forall \mathbf{x} \in \mathcal{X}, \forall \mathbf{y} \in \mathcal{Y}$

$$\mathcal{L}(\mathbf{x}^*, \mathbf{y}) \leq \mathcal{L}(\mathbf{x}^*, \mathbf{y}^*) \leq \mathcal{L}(\mathbf{x}, \mathbf{y}^*). \quad (2)$$

Examples of saddle point problems. Taskar et al. (2006) cast the maximum-margin estimation of structured output models as a bilinear saddle point problem $\mathcal{L}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top M \mathbf{y}$, where \mathcal{X} is the regularized set of parameters and \mathcal{Y} is an encoding of the set of possible structured outputs. They considered settings where projection on \mathcal{X} and \mathcal{Y} were efficient but one can imagine many situations where only LMO's are efficient. For example, we could use a structured sparsity inducing norm (Martins et al., 2011) for the parameter \mathbf{x} , such as the overlapping group lasso for which the projection is expensive (Bach et al., 2012), while \mathcal{Y} could be a combinatorial object such as the ground state of a planar Ising model (without external field) which admits an efficient oracle (Barahona, 1982) but has potentially intractable projection.

Similarly, two-player games (Von Neumann and Morgenstern, 1944) can often be solved as bilinear minimax problems. When a strategy space involves a polynomial number of constraints, the equilibria of such games can be solved efficiently (Koller et al., 1994). However, in situations such as the Colonel Blotto game or the Matching Duel (Ahmadinejad et al., 2016) the strategy space is intractably large and defined by an exponential number of linear constraints. Fortunately, some linear minimization oracles such as the blossom algorithm (Edmonds, 1965) can efficiently optimize over matching polytopes despite an exponential number of linear constraints.

Robust learning is also often cast as a saddle point minimax problem (Kim et al., 2005). Once again, a FW implementation could leverage fast linear oracles

while projection methods would be plagued by slower or intractable sub-problems. For instance, if the LMO is max-flow, it could have almost linear runtime while the corresponding projection would require cubic runtime quadratic programming (Kelner et al., 2014).

Related work. The standard approaches to solve smooth constrained saddle point problems are projection-type methods (surveyed in Xiu and Zhang (2003)), with in particular variations of Korpelevich’s extragradient method (Korpelevich, 1976), such as (Nesterov, 2007) which was used to solve the structured prediction problem (Taskar et al., 2006) mentioned above. There is surprisingly little work on FW-type methods for saddle point problems, although they were briefly considered for the more general *variational inequality* problem (VIP):

$$\text{find } \mathbf{z}^* \in \mathcal{Z} \text{ s.t. } \langle \mathbf{r}(\mathbf{z}^*), \mathbf{z} - \mathbf{z}^* \rangle \geq 0, \quad \forall \mathbf{z} \in \mathcal{Z}, \quad (3)$$

where \mathbf{r} is a Lipschitz mapping from \mathbb{R}^p to itself and $\mathcal{Z} \subseteq \mathbb{R}^p$. By using $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ and $\mathbf{r}(\mathbf{z}) = (\nabla_x \mathcal{L}(\mathbf{z}), -\nabla_y \mathcal{L}(\mathbf{z}))$, the variational inequality problem (3) reduces to the equivalent optimality conditions for the saddle point problem (1). Hammond (1984) showed that a FW algorithm with a step size of $O(1/t)$ converges for the VIP (3) when the set \mathcal{Z} is strongly convex, while FW with a generalized line-search on a saddle point problem is sometimes non-convergent when \mathcal{Z} is a polytope (see also (Patriksson, 1999, § 3.1.1)). She conjectured though that using a step size of $O(1/t)$ was also convergent when \mathcal{Z} is a polytope – a problem left open up to this point. More recently, Juditsky and Nemirovski (2016) (see also Cox et al. (2015)) proposed a method to transform a VIP on \mathcal{Z} where one has only access to a LMO, to a “dual” VIP on which they can use a projection-type method. Lan (2013) proposes to solve the saddle point problem (1) by running FW on \mathcal{X} on the *smoothed* version of the problem $\max_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}(\mathbf{x}, \mathbf{y})$, thus requiring a projection oracle on \mathcal{Y} . In contrast, in this paper we study simple approaches that do not require any transformations of the problem (1) nor any projection oracle on \mathcal{X} or \mathcal{Y} .

Contributions. In § 2, we extend several variants of the FW algorithm to solve the saddle point problem (1) that we think could be of interest to the machine learning community. In § 3, we give a first proof of (geometric) convergence for these methods over polytope domains under the assumptions of sufficient strong convex-concavity of \mathcal{L} , giving a partial answer to the conjecture from Hammond (1984). In § 4, we extend and refine the previous convergence results when \mathcal{X} and \mathcal{Y} are strongly convex sets and the gradient of \mathcal{L} is non-zero over $\mathcal{X} \times \mathcal{Y}$, while we survey the pure bilinear case in § 5. We finally present illustrative experiments

for our theory in § 6, noticing that the convergence theory is still incomplete for these methods.

2 Saddle point Frank-Wolfe (SP-FW)

The algorithms. This article will explore three SP extensions of the classical *Frank-Wolfe* algorithm (Algorithm 1) which are summarized in Algorithm 2, 3 and 4. In the following, the point computed by these algorithms after t steps will be noted $\mathbf{z}^{(t)} = (\mathbf{x}^{(t)}, \mathbf{y}^{(t)})$. We first obtain the *saddle point FW* (SP-FW) algorithm (Algorithm 2) by simultaneously doing a FW update on both convex functions $\mathcal{L}(\cdot, \mathbf{y}^{(t)})$ and $-\mathcal{L}(\mathbf{x}^{(t)}, \cdot)$ with a properly chosen step size. Hence the point $\mathbf{z}^{(t)}$ has a sparse representation as a convex combination of the points previously given by the oracle. This set of points is called the *active set*. If we assume that \mathcal{X} and \mathcal{Y} are the convex hulls of two finite sets of points \mathcal{A} and \mathcal{B} , we can also extend the *away-step Frank-Wolfe* (AFW) algorithm (Lacoste-Julien and Jaggi, 2015) to saddle point problems. As for AFW, this new algorithm is able to remove mass from “bad” atoms in the active set, (see L9 of Algorithm 3 and in Appendix A) to avoid the zig-zagging problem that slows down standard FW (Lacoste-Julien and Jaggi, 2015). Because of the special product structure of the domain, we consider more away directions than proposed in (Lacoste-Julien and Jaggi, 2015) for AFW. Namely, for every corner $\mathbf{v} = (\mathbf{v}_x, \mathbf{v}_y)$ and $\mathbf{v}' = (\mathbf{v}'_x, \mathbf{v}'_y)$ already picked, $\mathbf{x} - \mathbf{v}_x$ is a feasible directions in \mathcal{X} and $\mathbf{y} - \mathbf{v}'_y$ is a feasible direction in \mathcal{Y} . Thus the combination $(\mathbf{x} - \mathbf{v}_x, \mathbf{y} - \mathbf{v}'_y)$ is a feasible direction even if the particular corners \mathbf{v}_x and \mathbf{v}'_y have never been picked together. We thus maintain the iterates on \mathcal{X} and \mathcal{Y} as independent convex combination of their respective active sets of corners (Line 13 of Algorithm 3), i.e.,

$$\mathbf{x}^{(t)} = \sum_{\mathbf{v}_x \in \mathcal{S}_x^{(t)}} \alpha_{\mathbf{v}_x} \mathbf{v}_x \quad \text{and} \quad \mathbf{y}^{(t)} = \sum_{\mathbf{v}_y \in \mathcal{S}_y^{(t)}} \alpha_{\mathbf{v}_y} \mathbf{v}_y. \quad (4)$$

Finally, a straightforward saddle point generalization for the *pairwise Frank-Wolfe* (PFW) algorithm is given in Algorithm 4. The proposed algorithms preserve several nice properties of previous FW methods (in addition to only requiring LMO’s): simplicity of implementation, affine invariance (Jaggi, 2013), gap certificates computed for free, sparse representation of the iterates and the possibility to have adaptive step sizes using the gap computation. We next analyze the convergence of these algorithms.

The suboptimality error and the gap. To establish convergence, we first define several quantities of interest. In classical convex optimization, the suboptimality error h_t is well defined as $h_t := f(\mathbf{x}^{(t)}) - \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$. This quantity is clearly non-negative and

Algorithm 1 Frank-Wolfe algorithm

- 1: Let $\mathbf{x}^{(0)} \in \mathcal{X}$
- 2: **for** $t = 0 \dots T$ **do**
- 3: Compute $\mathbf{r}^{(t)} = \nabla f(\mathbf{x}^{(t)})$
- 4: Compute $\mathbf{s}^{(t)} := \operatorname{argmin}_{\mathbf{s} \in \mathcal{X}} \langle \mathbf{s}, \mathbf{r}^{(t)} \rangle$
- 5: Compute $g_t := \langle \mathbf{x}^{(t)} - \mathbf{s}^{(t)}, \mathbf{r}^{(t)} \rangle$
- 6: **if** $g_t \leq \epsilon$ **then return** $\mathbf{x}^{(t)}$
- 7: Let $\gamma = \frac{2}{2+t}$ (or do line-search)
- 8: Update $\mathbf{x}^{(t+1)} := (1 - \gamma)\mathbf{x}^{(t)} + \gamma\mathbf{s}^{(t)}$
- 9: **end for**

Algorithm 2 Saddle point Frank-Wolfe algorithm

- 1: Let $\mathbf{z}^{(0)} = (\mathbf{x}^{(0)}, \mathbf{y}^{(0)}) \in \mathcal{X} \times \mathcal{Y}$
- 2: **for** $t = 0 \dots T$ **do**
- 3: Compute $\mathbf{r}^{(t)} := \begin{pmatrix} \nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \\ -\nabla_y \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \end{pmatrix}$
- 4: Compute $\mathbf{s}^{(t)} := \operatorname{argmin}_{\mathbf{z} \in \mathcal{X} \times \mathcal{Y}} \langle \mathbf{z}, \mathbf{r}^{(t)} \rangle$
- 5: Compute $g_t := \langle \mathbf{z}^{(t)} - \mathbf{s}^{(t)}, \mathbf{r}^{(t)} \rangle$
- 6: **if** $g_t \leq \epsilon$ **then return** $\mathbf{z}^{(t)}$
- 7: Let $\gamma = \min\left(1, \frac{\nu}{2C} g_t\right)$ **or** $\gamma = \frac{2}{2+t}$
- 8: Update $\mathbf{z}^{(t+1)} := (1 - \gamma)\mathbf{z}^{(t)} + \gamma\mathbf{s}^{(t)}$
- 9: **end for**

Algorithm 3 Saddle point away-step Frank-Wolfe algorithm: **SP-AFW**($\mathbf{z}^{(0)}, \mathcal{A} \times \mathcal{B}, \epsilon$)

- 1: Let $\mathbf{z}^{(0)} = (\mathbf{x}^{(0)}, \mathbf{y}^{(0)}) \in \mathcal{A} \times \mathcal{B}$, $\mathcal{S}_x^{(0)} := \{\mathbf{x}^{(0)}\}$ and $\mathcal{S}_y^{(0)} := \{\mathbf{y}^{(0)}\}$
- 2: **for** $t = 0 \dots T$ **do**
- 3: Let $\mathbf{s}^{(t)} := \operatorname{LMO}_{\mathcal{A} \times \mathcal{B}}(\mathbf{r}^{(t)})$ and $\mathbf{d}_{\text{FW}}^{(t)} := \mathbf{s}^{(t)} - \mathbf{z}^{(t)}$ ($\mathbf{r}^{(t)}$ as defined in L3 in Algorithm 2)
- 4: Let $\mathbf{v}^{(t)} \in \operatorname{argmax}_{\mathbf{v} \in \mathcal{S}_x^{(t)} \times \mathcal{S}_y^{(t)}} \langle \mathbf{r}^{(t)}, \mathbf{v} \rangle$ and $\mathbf{d}_A^{(t)} := \mathbf{z}^{(t)} - \mathbf{v}^{(t)}$ (the away direction)
- 5: **if** $g_t^{\text{FW}} := \langle -\mathbf{r}^{(t)}, \mathbf{d}_{\text{FW}}^{(t)} \rangle \leq \epsilon$ **then return** $\mathbf{z}^{(t)}$ (FW gap is small enough, so return)
- 6: **if** $\langle -\mathbf{r}^{(t)}, \mathbf{d}_{\text{FW}}^{(t)} \rangle \geq \langle -\mathbf{r}^{(t)}, \mathbf{d}_A^{(t)} \rangle$ **then**
- 7: $\mathbf{d}^{(t)} := \mathbf{d}_{\text{FW}}^{(t)}$, and $\gamma_{\max} := 1$ (choose the FW direction)
- 8: **else**
- 9: $\mathbf{d}^{(t)} := \mathbf{d}_A^{(t)}$, and $\gamma_{\max} := \min \left\{ \frac{\alpha_{\mathbf{v}_x^{(t)}}}{1 - \alpha_{\mathbf{v}_x^{(t)}}}, \frac{\alpha_{\mathbf{v}_y^{(t)}}}{1 - \alpha_{\mathbf{v}_y^{(t)}}} \right\}$ (maximum feasible step size; a drop step is when $\gamma_t = \gamma_{\max}$)
- 10: **end if**
- 11: Let $g_t^{\text{PFW}} = \langle -\mathbf{r}^{(t)}, \mathbf{d}_{\text{FW}}^{(t)} + \mathbf{d}_A^{(t)} \rangle$ **and** $\gamma_t = \min \left\{ \gamma_{\max}, \frac{\nu^{\text{PFW}}}{2C} g_t^{\text{PFW}} \right\}$ (ν^{PFW} and C set as in Thm. 1)
- 12: Update $\mathbf{z}^{(t+1)} := \mathbf{z}^{(t)} + \gamma_t \mathbf{d}^{(t)}$ (and accordingly for the weights $\alpha^{(t+1)}$, see Lacoste-Julien and Jaggi (2015))
- 13: Update $\mathcal{S}_x^{(t+1)} := \{\mathbf{v}_x \in \mathcal{A} \text{ s.t. } \alpha_{\mathbf{v}_x}^{(t+1)} > 0\}$ and $\mathcal{S}_y^{(t+1)} := \{\mathbf{v}_y \in \mathcal{B} \text{ s.t. } \alpha_{\mathbf{v}_y}^{(t+1)} > 0\}$
- 14: **end for**

Algorithm 4 Saddle point pairwise Frank-Wolfe algorithm: **SP-PFW**($\mathbf{z}^{(0)}, \mathcal{A} \times \mathcal{B}, \epsilon$)

- 1: In Alg. 3, replace L6 to L10 by: $\mathbf{d}^{(t)} := \mathbf{d}_{\text{PFW}}^{(t)} := \mathbf{s}^{(t)} - \mathbf{v}^{(t)}$, and $\gamma_{\max} := \min \left\{ \alpha_{\mathbf{v}_x^{(t)}}, \alpha_{\mathbf{v}_x^{(t)}} \right\}$.

proving that h_t goes to 0 is enough to establish convergence. Unfortunately, in the saddle point setting the quantity $\mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}^*$ is no longer non-negative and can be equal to zero for an infinite number of points (\mathbf{x}, \mathbf{y}) while $(\mathbf{x}, \mathbf{y}) \notin (\mathcal{X}^*, \mathcal{Y}^*)$. For instance, if $\mathcal{L}(\mathbf{x}, \mathbf{y}) = \mathbf{x} \cdot \mathbf{y}$ with $\mathcal{X} = \mathcal{Y} = [-1, 1]$, then $\mathcal{L}^* = 0$ and $(\mathcal{X}^*, \mathcal{Y}^*) = \{(0, 0)\}$. But for all $\mathbf{x} \in \mathcal{X}$ and $\mathbf{y} \in \mathcal{Y}$, $\mathbf{x} \cdot 0 = 0 \cdot \mathbf{y} = \mathcal{L}^*$. The saddle point literature thus considers a non-negative gap function (also known as a merit function (Larsson and Patriksson, 1994; Zhu and Marcotte, 1998) and (Patriksson, 1999, Sec 4.4.1)) which is zero only for optimal points, in order to quantify progress towards the saddle point. We can define the following *suboptimality error* h_t for our saddle point

problem:

$$\begin{aligned}
 h_t &:= \mathcal{L}(\mathbf{x}^{(t)}, \hat{\mathbf{y}}^{(t)}) - \mathcal{L}(\hat{\mathbf{x}}^{(t)}, \mathbf{y}^{(t)}), \\
 \text{where } \hat{\mathbf{x}}^{(t)} &:= \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \mathcal{L}(\mathbf{x}, \mathbf{y}^{(t)}), \\
 \text{and } \hat{\mathbf{y}}^{(t)} &:= \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}).
 \end{aligned} \tag{5}$$

This is an example of *primal-dual* gap function by noticing that

$$\begin{aligned}
 h_t &= \mathcal{L}(\hat{\mathbf{y}}^{(t)}, \mathbf{x}^{(t)}) - \mathcal{L}^* + \mathcal{L}^* - \mathcal{L}(\mathbf{y}^{(t)}, \hat{\mathbf{x}}^{(t)}) \\
 &= p(\mathbf{x}^{(t)}) - p(\mathbf{x}^*) + g(\mathbf{y}^*) - g(\mathbf{y}^{(t)}),
 \end{aligned} \tag{6}$$

where $p(\mathbf{x}) := \max_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}(\mathbf{x}, \mathbf{y})$ is the convex primal function and $g(\mathbf{y}) := \min_{\mathbf{x} \in \mathcal{X}} \mathcal{L}(\mathbf{x}, \mathbf{y})$ is the concave dual function. By convex-concavity, h_t can be upper-bounded by the following FW linearization gap (Jaggi,

2011, 2013; Larsson and Patriksson, 1994; Zhu and Marcotte, 1998):

$$g_t^{\text{FW}} := \max_{\mathbf{s}_x \in \mathcal{X}} \left\langle \mathbf{x}^{(t)} - \mathbf{s}_x, \nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \right\rangle \Big\} := g_t^{(x)} \\ + \max_{\mathbf{s}_y \in \mathcal{Y}} \left\langle \mathbf{y}^{(t)} - \mathbf{s}_y, -\nabla_y \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \right\rangle \Big\} := g_t^{(y)}. \quad (7)$$

This gap is easy to compute and gives a stopping criterion since $g_t^{\text{FW}} \geq h_t$.

Compensation phenomenon and difficulty for SP. Even when equipped with a suboptimality error and a gap function (as in the convex case), we still cannot apply the standard FW convergence analysis. The usual FW proof sketch uses the fact that the gradient of f is Lipschitz continuous to get

$$h_{t+1} \leq h_t - \gamma_t g_t^{\text{FW}} + \gamma_t^2 \frac{L \|\mathbf{d}^{(t)}\|^2}{2} \quad (8)$$

which then provides a rate of convergence. Roughly, since $g_t \geq h_t$ by convexity, if γ_t is small enough then (h_t) will decrease and converge. For simplicity, in the main paper, $\|\cdot\|$ will refer to the ℓ_2 norm of \mathbb{R}^d . The partial Lipschitz constants and the diameters of the sets are defined with respect to this norm (see (40) in Appendix B.1 for more general norms).

Using the L -Lipschitz continuity of \mathcal{L} and letting $\mathcal{L}_t := \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})$ as a shorthand, we get

$$\mathcal{L}_{t+1} \leq \mathcal{L}_t + \gamma_t \left\langle \mathbf{d}_x^{(t)}, \nabla_x \mathcal{L}_t \right\rangle + \gamma_t \left\langle \mathbf{d}_y^{(t)}, \nabla_y \mathcal{L}_t \right\rangle \\ + \gamma_t^2 \frac{L \|\mathbf{d}^{(t)}\|^2}{2} \quad (9)$$

where $\mathbf{d}_x^{(t)} = \mathbf{s}_x^{(t)} - \mathbf{x}^{(t)}$ and $\mathbf{d}_y^{(t)} = \mathbf{s}_y^{(t)} - \mathbf{y}^{(t)}$. Then

$$\mathcal{L}_{t+1} - \mathcal{L}^* \leq \mathcal{L}_t - \mathcal{L}^* - \gamma_t (g_t^{(x)} - g_t^{(y)}) + \gamma_t^2 \frac{L \|\mathbf{d}^{(t)}\|^2}{2}. \quad (10)$$

Unfortunately, the quantity g_t^{FW} does *not* appear above and we therefore cannot control the oscillation of the sequence (the quantity $g_t^{(x)} - g_t^{(y)}$ can make the sequence increase or decrease). Instead, we must focus on more specific SP optimization settings and introduce other quantities of interest in order to establish convergence.

The asymmetry of the SP. Hammond (1984, p. 165) showed the divergence of the SP-FW algorithm with an extended line-search step-size on some bilinear objectives. She mentioned that the difficulty for SP optimization is contained in this bilinear coupling between \mathbf{x} and \mathbf{y} . More generally, most of the examples of SP functions cited in the introduction can be written in the form:

$$\mathcal{L}(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{x}^\top M \mathbf{y} - g(\mathbf{y}), \quad f \text{ and } g \text{ convex.} \quad (11)$$

In this setting, the bilinear part M is the only term preventing us to apply theorems on standard FW. Hammond (1984, p. 175) also conjectured that the SP-FW algorithm with $\gamma_t = 1/(t+1)$ performed on a uniformly strongly convex-concave objective function (see (12)) over a polytope should converge. We give a partial answer to this conjecture in the following section.

3 SP-FW for strongly convex functions

Uniform strong convex-concavity. In this section, we will assume that \mathcal{L} is uniformly $(\mu_{\mathcal{X}}, \mu_{\mathcal{Y}})$ -strongly convex-concave, which means that the following function is convex-concave:

$$(\mathbf{x}, \mathbf{y}) \mapsto \mathcal{L}(\mathbf{x}, \mathbf{y}) - \frac{\mu_{\mathcal{X}}}{2} \|\mathbf{x}\|^2 + \frac{\mu_{\mathcal{Y}}}{2} \|\mathbf{y}\|^2. \quad (12)$$

A new merit function. To prove our theorem, we use a different quantity w_t which is smaller than h_t but still a valid merit function in the case of *strongly convex-concave* SPs (where $(\mathbf{x}^*, \mathbf{y}^*)$ is thus unique); see (14) below. For $(\mathbf{x}^*, \mathbf{y}^*)$ a solution of (1), we define the non-negative quantity w_t :

$$w_t := \underbrace{\mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^*) - \mathcal{L}^*}_{:=w_t^{(x)}} + \underbrace{\mathcal{L}^* - \mathcal{L}(\mathbf{x}^*, \mathbf{y}^{(t)})}_{:=w_t^{(y)}}. \quad (13)$$

Notice that $w_t^{(x)}$ and $w_t^{(y)}$ are non-negative, and that $w_t \leq h_t$ since:

$$\mathcal{L}(\mathbf{x}^{(t)}, \hat{\mathbf{y}}^{(t)}) - \mathcal{L}(\hat{\mathbf{x}}^{(t)}, \mathbf{y}^{(t)}) \geq \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^*) - \mathcal{L}(\mathbf{x}^*, \mathbf{y}^{(t)}).$$

In general, w_t can be zero even if we have not reached a solution. For example, with $\mathcal{L}(\mathbf{x}, \mathbf{y}) = \mathbf{x} \cdot \mathbf{y}$ and $\mathcal{X} = \mathcal{Y} = [-1, 1]$, then $\mathbf{x}^* = \mathbf{y}^* = \mathbf{0}$, implying $w_t = 0$ for any $(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})$. But for a uniformly strongly convex-concave \mathcal{L} , this cannot happen and we can prove that w_t has the following nice property (Proposition 15 in Appendix B.6):

$$h_t \leq \sqrt{2} P_{\mathcal{L}} \sqrt{w_t}, \quad (14)$$

where

$$P_{\mathcal{L}} \leq \sqrt{2} \sup_{\mathbf{z} \in \mathcal{X} \times \mathcal{Y}} \left\{ \frac{\|\nabla_x \mathcal{L}(\mathbf{z})\|_{\mathcal{X}^*}}{\sqrt{\mu_{\mathcal{X}}}}, \frac{\|\nabla_y \mathcal{L}(\mathbf{z})\|_{\mathcal{Y}^*}}{\sqrt{\mu_{\mathcal{Y}}}} \right\}. \quad (15)$$

Pyramidal width and distance to the border. We now provide a theorem that establishes convergence in two situations: **(I)** when the SP belongs to the interior of $\mathcal{X} \times \mathcal{Y}$; **(P)** when the set is a polytope, i.e. when there exist two finite sets such that $\mathcal{X} = \text{conv}(\mathcal{A})$ and $\mathcal{Y} = \text{conv}(\mathcal{B})$. Our convergence result holds when (roughly) the strong convex-concavity of \mathcal{L} is big enough

in comparison to the cross Lipschitz constants L_{XY} , L_{YX} of $\nabla\mathcal{L}$ (defined in (20) below) multiplied by geometric “condition numbers” of each set. The condition number of \mathcal{X} (and similarly for \mathcal{Y}) is defined as the ratio of its diameter $D_{\mathcal{X}} := \sup_{\mathbf{x}, \mathbf{x}' \in \mathcal{X}} \|\mathbf{x} - \mathbf{x}'\|$ over the following appropriate notions of “width”:

$$\text{border distance: } \delta_{\mathcal{X}} := \min_{\mathbf{s} \in \partial\mathcal{X}} \|\mathbf{x}^* - \mathbf{s}\| \text{ for (I), (16)}$$

$$\text{pyramidal width: } \delta_{\mathcal{A}} := PWidth(\mathcal{A}) \text{ for (P). (17)}$$

The pyramidal width (17) is formally defined in Eq. 9 of Lacoste-Julien and Jaggi (2015). Given the above constants, we can state below a non-affine invariant version of our convergence theorem (for simplicity). The affine invariant versions of this theorem are given in Thm. 24 and 25 in Appendix D.2 (with proofs).

Theorem 1. *Let \mathcal{L} be a convex-concave function and $\mathcal{X} \times \mathcal{Y}$ a convex and compact set. Assume that the gradient of \mathcal{L} is L -Lipschitz continuous, that \mathcal{L} is $(\mu_{\mathcal{X}}, \mu_{\mathcal{Y}})$ -strongly convex-concave, and that we are in one of the two following situations:*

(I) *The SP belongs to the interior of $\mathcal{X} \times \mathcal{Y}$. In this case, set $g_t = g_t^{\text{FW}}$ (as in L5 of Alg. 3), $\delta_{\mu} := \sqrt{\min(\mu_{\mathcal{X}}\delta_{\mathcal{X}}^2, \mu_{\mathcal{Y}}\delta_{\mathcal{Y}}^2)}$ and $a := 1$. “Algorithm” then refers to SP-FW.*

(P) *The sets \mathcal{X} and \mathcal{Y} are polytopes. In this case, set $g_t = g_t^{\text{PFW}}$ (as in L11 of Alg. 3), $\delta_{\mu} := \sqrt{\min(\mu_{\mathcal{X}}\delta_{\mathcal{A}}^2, \mu_{\mathcal{Y}}\delta_{\mathcal{B}}^2)}$ and $a := \frac{1}{2}$. “Algorithm” then refers to SP-AFW. Here δ_{μ} needs to use the Euclidean norm for its defining constants.*

In both cases, if $\nu := a - \frac{\sqrt{2}}{\delta_{\mu}} \max\left\{\frac{D_{\mathcal{X}}L_{XY}}{\sqrt{\mu_{\mathcal{Y}}}}, \frac{D_{\mathcal{Y}}L_{YX}}{\sqrt{\mu_{\mathcal{X}}}}\right\}$ is positive, then the errors h_t (5) of the iterates of the algorithm with step size $\gamma_t = \min\{\gamma_{\max}, \frac{\nu}{2C}g_t\}$ decrease geometrically as

$$h_t = O\left((1 - \rho)^{\frac{k(t)}{2}}\right) \text{ and } \min_{s \leq t} g_s = O\left((1 - \rho)^{\frac{k(t)}{2}}\right)$$

where $\rho := \frac{\nu^2\mu}{2C}$, $C := \frac{LD_{\mathcal{X}}^2 + LD_{\mathcal{Y}}^2}{2}$ and $k(t)$ is the number of non-drop step after t steps (see L9 in Alg. 3). In case (I) we have $k(t) = t$ and in case (P) we have $k(t) \geq t/3$. For both algorithms, if $\delta_{\mu} > 2 \max\left\{\frac{D_{\mathcal{X}}L_{XY}}{\mu_{\mathcal{X}}}, \frac{D_{\mathcal{Y}}L_{YX}}{\mu_{\mathcal{Y}}}\right\}$, we also obtain a sublinear rate with the universal choice $\gamma_t = \min\{\gamma_{\max}, \frac{2}{2+k(t)}\}$. This yields the rates:

$$\min_{s \leq t} h_s \leq \min_{s \leq t} g_s^{\text{FW}} = O\left(\frac{1}{t}\right). \quad (18)$$

Clearly, the sublinear rate seems less interesting than the linear one but has the added convenience that the step size can be set without knowledge of various constants that characterize \mathcal{L} . Moreover, it provides a partial answer to the conjecture from Hammond (1984).

Proof sketch. Strong convexity is an essential assumption in our proof; it allows us to relate w_t to how close we are to the optimum. Actually, by $\mu_{\mathcal{Y}}$ -strong concavity of $\mathcal{L}(\mathbf{x}^*, \cdot)$ we have

$$\|\mathbf{y}^{(t)} - \mathbf{y}^*\| \leq \sqrt{\frac{2}{\mu_{\mathcal{Y}}}(\mathcal{L}^* - \mathcal{L}(\mathbf{x}^*, \mathbf{y}^{(t)}))} = \sqrt{\frac{2}{\mu_{\mathcal{Y}}}w_t^{(y)}}. \quad (19)$$

Now, recall that we assumed that $\nabla\mathcal{L}$ is Lipschitz continuous. In the following, we will call L the Lipschitz continuity constant of $\nabla\mathcal{L}$ and L_{XY} and L_{YX} its (cross) partial Lipschitz constants. For all $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$, $\mathbf{y}, \mathbf{y}' \in \mathcal{Y}$, these constants satisfy

$$\begin{aligned} \|\nabla_{\mathbf{x}}\mathcal{L}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{x}}\mathcal{L}(\mathbf{x}, \mathbf{y}')\|_{\mathcal{X}^*} &\leq L_{XY}\|\mathbf{y} - \mathbf{y}'\|_{\mathcal{Y}}, \\ \|\nabla_{\mathbf{y}}\mathcal{L}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{y}}\mathcal{L}(\mathbf{x}', \mathbf{y})\|_{\mathcal{Y}^*} &\leq L_{YX}\|\mathbf{x} - \mathbf{x}'\|_{\mathcal{X}}. \end{aligned} \quad (20)$$

Note that $L_{XY}, L_{YX} \leq L$ if $\|(\mathbf{x}, \mathbf{y})\| := \|\mathbf{x}\|_{\mathcal{X}} + \|\mathbf{y}\|_{\mathcal{Y}}$. Then, using Lipschitz continuity of the gradient,

$$\begin{aligned} \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^*) &\leq \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^*) + \gamma\langle \mathbf{d}_x^{(t)}, \nabla_{\mathbf{x}}\mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^*) \rangle \\ &\quad + \gamma^2 \frac{L\|\mathbf{d}_x^{(t)}\|^2}{2}. \end{aligned} \quad (21)$$

Furthermore, setting $(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^{(t)}, \mathbf{y}^*)$ and $\mathbf{y}' = \mathbf{y}^{(t)}$ in Equation (20), we have

$$\begin{aligned} w_{t+1}^{(x)} &\leq w_t^{(x)} - \gamma g_t^{(x)} + \gamma D_{\mathcal{X}}L_{XY}\|\mathbf{y}^{(t)} - \mathbf{y}^*\| \\ &\quad + \gamma^2 \frac{LD_{\mathcal{X}}^2}{2}. \end{aligned} \quad (22)$$

Finally, combining (22) and (19), we get

$$\begin{aligned} w_{t+1}^{(x)} &\leq w_t^{(x)} - \gamma g_t^{(x)} + \gamma D_{\mathcal{X}}L_{XY}\sqrt{\frac{2}{\mu_{\mathcal{Y}}}}\sqrt{w_t^{(y)}} \\ &\quad + \gamma^2 \frac{LD_{\mathcal{X}}^2}{2}. \end{aligned} \quad (23)$$

A similar argument on $-\mathcal{L}(\mathbf{x}^*, \mathbf{y}^{(t+1)})$ gives a bound on $w_t^{(y)}$ much like (23). Summing both yields:

$$\begin{aligned} w_{t+1} &\leq w_t - \gamma g_t + 2\gamma \max\left\{\frac{D_{\mathcal{X}}L_{XY}}{\sqrt{\mu_{\mathcal{Y}}}}, \frac{D_{\mathcal{Y}}L_{YX}}{\sqrt{\mu_{\mathcal{X}}}}\right\}\sqrt{w_t} \\ &\quad + \gamma^2 \frac{LD_{\mathcal{X}}^2 + LD_{\mathcal{Y}}^2}{2}. \end{aligned} \quad (24)$$

We now apply recent developments in the convergence theory of FW methods for strongly convex objectives. Lacoste-Julien and Jaggi (2015) crucially upper bound the square root of the suboptimality error on a convex function with the FW gap if the optimum is in the interior, or with the PFW gap if the set is a polytope (Lemma 18 in Appendix C.2). We continue our proof sketch for case (I) only¹:

$$2\mu_{\mathcal{X}}\delta_{\mathcal{X}}^2 \left(\mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^*, \mathbf{y}^{(t)})\right) \leq \left(g_t^{(x)}\right)^2 \quad (25)$$

$$\text{where } \delta_{\mathcal{X}} := \min_{\mathbf{x} \in \partial\mathcal{X}} \|\mathbf{x}^* - \mathbf{s}\|.$$

¹The idea is similar for case (P), but with the additional complication of possible drop steps.

We can also get the respective equation on \mathbf{y} with $\delta_{\mathcal{Y}} := \min_{\mathbf{y} \in \partial \mathcal{Y}} \|\mathbf{y}^* - \mathbf{y}\|$ and sum it with the previous one (25) to get:

$$\delta_{\mu} \sqrt{2w_t} \leq g_t \quad \text{where } \delta_{\mu} := \sqrt{\min(\mu_{\mathcal{X}} \delta_{\mathcal{X}}^2, \mu_{\mathcal{Y}} \delta_{\mathcal{Y}}^2)}. \quad (26)$$

Plugging this last equation into (23) gives us

$$w_{t+1} \leq w_t - \nu \gamma g_t + \gamma^2 C \quad \text{where } C := \frac{LD_{\mathcal{X}}^2 + LD_{\mathcal{Y}}^2}{2} \quad (27)$$

and $\nu := 1 - \frac{\sqrt{2}}{\delta_{\mu}} \max \left\{ \frac{D_{\mathcal{X}} L_{\mathcal{X}\mathcal{Y}}}{\sqrt{\mu_{\mathcal{Y}}}}, \frac{D_{\mathcal{Y}} L_{\mathcal{Y}\mathcal{X}}}{\sqrt{\mu_{\mathcal{X}}}} \right\}$.

The recurrence (27) is typical in the FW literature. We can re-apply standard techniques on the sequence w_t to get a sublinear rate with $\gamma_t = \frac{2}{2+t}$, or a linear rate with $\gamma_t = \min \left\{ \gamma_{\max}, \frac{\nu g_t}{2C} \right\}$ (which minimizes the RHS of (27) and actually guarantees that w_t will be decreasing). Finally, thanks to strong convexity, a rate on w_t gives us a rate on h_t (by (14)). \square

4 SP-FW with strongly convex sets

Strongly convex set. One can (roughly) define strongly convex set as sublevel sets of strongly convex functions (Vial, 1983, Prop. 4.14). In this section, we replace the strong convex-concavity assumption on \mathcal{L} with the assumption that \mathcal{X} and \mathcal{Y} are β -strongly convex sets.

Definition 2 (Vial (1983); Polyak (1966)). *A convex set \mathcal{X} is said to be β -strongly convex with respect to $\|\cdot\|$ if for any $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ and any $\gamma \in [0, 1]$, $B_{\beta}(\gamma, \mathbf{x}, \mathbf{y}) \subset \mathcal{X}$ where $B_{\beta}(\gamma, \mathbf{x}, \mathbf{y})$ is the $\|\cdot\|$ -ball of radius $\gamma(1-\gamma)\frac{\beta}{2}\|\mathbf{x}-\mathbf{y}\|^2$ centered at $\gamma\mathbf{x} + (1-\gamma)\mathbf{y}$.*

Frank-Wolfe for convex optimization over strongly convex sets has been studied by Levitin and Polyak (1966); Demyanov and Rubinov (1970) and Dunn (1979), amongst others. They all obtained a linear rate for the FW algorithm if the norm of the gradient is lower bounded by a constant. More recently, Garber and Hazan (2015) proved a sublinear rate $O(1/t^2)$ by replacing the lower bound on the gradient by a strong convexity assumption on the function. In the VIP setting (3), the linear convergence has been proved if the optimization is done under a strongly convex set but this assumption does *not* extend to $\mathcal{X} \times \mathcal{Y}$ which *cannot* be strongly convex if \mathcal{X} or \mathcal{Y} is not reduced to a single element. In order to prove the convergence, we first prove the Lipschitz continuity of the *FW-corner* function $\mathbf{s}(\cdot)$ defined below. A proof of this theorem is given in Appendix E.

Theorem 3. *Let \mathcal{X} and \mathcal{Y} be β -strongly convex sets. If $\min(\|\nabla_{\mathbf{x}} L(\mathbf{z})\|_{\mathcal{X}^*}, \|\nabla_{\mathbf{y}} L(\mathbf{z})\|_{\mathcal{Y}^*}) \geq \delta > 0$ for all $\mathbf{z} \in \mathcal{X} \times \mathcal{Y}$, then the oracle function $\mathbf{z} \mapsto \mathbf{s}(\mathbf{z}) := \arg \min_{\mathbf{s} \in \mathcal{X} \times \mathcal{Y}} \langle \mathbf{s}, \mathbf{r}(\mathbf{z}) \rangle$ is well defined and is $\frac{4L}{\delta\beta}$ -Lipschitz continuous (using the norm $\|(\mathbf{x}, \mathbf{y})\|_{\mathcal{X} \times \mathcal{Y}} := \|\mathbf{x}\|_{\mathcal{X}} + \|\mathbf{y}\|_{\mathcal{Y}}$), where $\mathbf{r}(\mathbf{z}) := (\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{z}), -\nabla_{\mathbf{y}} \mathcal{L}(\mathbf{z}))$.*

Convergence rate. When the FW-corner function $\mathbf{s}(\cdot)$ is Lipschitz continuous (by Theorem 3), we can actually show that the FW gap is decreasing in the FW direction and get a similar inequality as the standard FW one (8), but, in this case, on the *gaps*: $g_{t+1} \leq g_t(1 - \gamma_t) + \gamma_t^2 \|\mathbf{s}^{(t)} - \mathbf{z}^{(t)}\|^2 C_{\delta}$. Moreover, one can show that the FW gap on a strongly convex set \mathcal{X} can be lower-bounded by $\|\mathbf{s}_x^{(t)} - \mathbf{x}^{(t)}\|^2$ (Lemma 27 in Appendix E), by using the fact that \mathcal{X} contains a ball of sufficient radius around the midpoint between $\mathbf{s}_x^{(t)}$ and $\mathbf{x}^{(t)}$. From these two facts, we can prove the following linear rate of convergence (*not* requiring any strong convex-concavity of \mathcal{L}).

Theorem 4. *Let \mathcal{L} be a convex-concave function and \mathcal{X} and \mathcal{Y} two compact β -strongly convex sets. Assume that the gradient of \mathcal{L} is L -Lipschitz continuous and that there exists $\delta > 0$ such that $\min(\|\nabla_{\mathbf{x}} L(\mathbf{z})\|_*, \|\nabla_{\mathbf{y}} L(\mathbf{z})\|_*) \geq \delta \forall \mathbf{z} \in \mathcal{X} \times \mathcal{Y}$. Set $C_{\delta} := 2L + \frac{8L^2}{\beta\delta}$. Then the gap g_t^{FW} of the SP-FW algorithm with step size $\gamma_t = \frac{g_t^{\text{FW}}}{\|\mathbf{s}^{(t)} - \mathbf{z}^{(t)}\|^2 C_{\delta}}$ converges linearly as $g_t^{\text{FW}} \leq g_0(1 - \rho)^t$, where $\rho := \frac{\beta\delta}{16C_{\delta}}$.*

5 SP-FW in the bilinear setting

Fictitious play. In her thesis, Hammond (1984, § 4.3.1) pointed out that for the bilinear setting:

$$\min_{\mathbf{x} \in \Delta_p} \max_{\mathbf{y} \in \Delta_q} \mathbf{x}^{\top} M \mathbf{y} \quad (28)$$

where Δ_p is the probability simplex on p elements, the SP-FW algorithm with step size $\gamma_t = 1/(1+t)$ is equivalent to the fictitious play (FP) algorithm introduced by Brown (1951). The FP algorithm has been widely studied in the game literature. Its convergence has been proved by Robinson (1951), while Shapiro (1958) showed that one can deduce from Robinson's proof a $O(t^{-1/(p+q-2)})$ rate. Around the same time, Karlin (1960) conjectured that the FP algorithm converged at the better rate of $O(t^{-1/2})$, though this conjecture is still open and Shapiro's rate is the only one we are aware of. Interestingly, Daskalakis and Pan (2014) recently showed that Shapiro's rate is also a lower bound if the tie breaking rule gets the worst pick an infinite number of times. Nevertheless, this kind of adversarial tie breaking rule does not seem realistic since this rule is a priori defined by the programmer. In practical cases (by setting a fixed prior order for ties or picking randomly for example), Karlin's Conjecture (Karlin, 1960) is still open. Moreover, we always observed an empirical rate of at least $O(t^{-1/2})$ during our experiments, we thus believe the conjecture to be true for realistic tie breaking rules.

Rate for SP-FW. Via the affine invariance of the FW algorithm and the fact that every polytope with

p vertices is the affine transformation of a probability simplex of dimension p , any rate for the fictitious play algorithm implies a rate for SP-FW.

Corollary 5. *For polytopes \mathcal{X} and \mathcal{Y} with p and q vertices respectively and $\mathcal{L}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top M \mathbf{y}$, the SP-FW algorithm with step size $\gamma_t = \frac{1}{t+1}$ converges at the rate $h_t = O\left(t^{-\frac{1}{p+q-2}}\right)$.*

This (very slow) convergence rate is mainly of theoretical interest, providing a safety check that the algorithm actually converges. Moreover, if Karlin’s strong conjecture is true, we can get a $O(1/\sqrt{t})$ worst case rate which is confirmed by our experiments.

6 Experiments

Toy experiments. First, we test the empirical convergence of our algorithms on a simple saddle point problem over the unit cube in dimension d (whose pyramidal width has the explicit value $1/\sqrt{d}$ by Lemma 4 from Lacoste-Julien and Jaggi (2015)). Thus $\mathcal{X} = \mathcal{Y} := [0, 1]^d$ and the linear minimization oracle is simply $\text{LMO}(\cdot) = -0.5 \cdot (\text{sign}(\cdot) - \mathbf{1})$. We consider the following objective function:

$$\frac{\mu}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2 + (\mathbf{x} - \mathbf{x}^*)^\top M (\mathbf{y} - \mathbf{y}^*) - \frac{\mu}{2} \|\mathbf{y} - \mathbf{y}^*\|_2^2 \quad (29)$$

for which we can control the location of the saddle point $(\mathbf{x}^*, \mathbf{y}^*) \in \mathcal{X} \times \mathcal{Y}$. We generate a matrix M randomly as $M \sim \mathcal{U}([-0.1, 0.1]^{d \times d})$ and keep it fixed for all experiments. For the interior point setup (I), we set $(\mathbf{x}^*, \mathbf{y}^*) \sim \mathcal{U}([0.25, 0.75]^{2d})$, while we set \mathbf{x}^* and \mathbf{y}^* to some fixed random vertex of the unit cube for the setup (P). With all these parameters fixed, the constant ν is a function of μ only. We thus vary the strong convexity parameter μ to test various ν ’s.

We verify the linear convergence expected for the SP-FW algorithm for case (I) in Figure 1a, and for the SP-AFW algorithm for case (P) in Figure 1b. As the adaptive step size (and rate) depends linearly on ν , the linear rate becomes quite slow for small ν . In this regime (in red), the step size $2/(2+k(t))$ (in orange) can actually perform better, despite its theoretical sublinear rate.

Finally, figure 1c shows that we can observe a linear convergence of SP-AFW even if ν is negative by using a different step size. In this case, we use the heuristic adaptive step size $\gamma_t := g_t/\tilde{C}$ where $\tilde{C} := LD_{\mathcal{X}}^2 + LD_{\mathcal{Y}}^2 + L_{XY}L_{YX} (D_{\mathcal{X}}^2/\mu_{\mathcal{X}} + D_{\mathcal{Y}}^2/\mu_{\mathcal{Y}})$. Here \tilde{C} takes into account the coupling between the concave and the convex variable and is motivated from a different proof of convergence that we were not able to complete. The empirical linear convergence in this case is not

yet supported by a complete analysis, highlighting the need for more sophisticated arguments.

Graphical games. We now consider a bilinear objective $\mathcal{L}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top M \mathbf{y}$ where exact projections on the sets is intractable, but we have a tractable LMO. The problem is motivated from the following setup. We consider a game between two universities (A and B) that are admitting s students and have to assign pairs of students into dorms. If students are unhappy with their dorm assignments, they will go to the other university. The game has a payoff matrix M belonging to $\mathbb{R}^{(s(s-1)/2)^2}$ where $M_{ij,kl}$ is the expected tuition that B gets (or A gives up) if A pairs student i with j and B pairs student k with l . Here the actions \mathbf{x} and \mathbf{y} are both in the marginal polytope of all perfect unipartite matchings. Assume that we are given a graph $G = (V, E)$ with vertices V and edges E . For a subset of nodes $S \subseteq V$, let the induced subgraph $G(S) = (S, E(S))$. Edmonds (1965) showed that any subgraph forming a triangle can contain at most one edge of any perfect matching. This forms an exponential set of linear equalities which define the matching polytope $\mathcal{P}(G) \subset \mathbb{R}^E$ as

$$\{\mathbf{x} \mid \mathbf{x}_e \geq 0, \sum_{e \in E(S)} \mathbf{x}_e \leq k, \forall S \subseteq V, |S| = 2k+1, \forall e \in E\}. \quad (30)$$

While this strategy space seems daunting, the LMO can be solved in $\mathcal{O}(s^3)$ time using the blossom algorithm (Edmonds, 1965). We run the SP-FW algorithm with $\gamma_t = 2/(t+2)$ on this problem with $s = 2^j$ students for $j = 3, \dots, 8$ with results given in Figure 1d ($d = s(s-1)/2$ in the legend represents the dimensionality of the \mathbf{x} and \mathbf{y} variables). The order of the complexity of the LMO is then $\mathcal{O}(d^{3/2})$. In Figure 1d, the observed empirical rate of the SP-FW algorithm (using $\gamma_t = 2/(t+2)$) is $\mathcal{O}(1/t^2)$. Empirically, faster rates seem to arise if the solution is at a corner (a pure equilibrium, to be expected for random payoff matrices in light of (Bárány et al., 2007)).

Sparse structured SVM. We finally consider a challenging optimization problem arising from structured prediction. We consider the saddle point formulation (Taskar et al., 2006) for a ℓ_1 -regularized structured SVM objective that minimizes the primal cost function $p(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \tilde{H}_i(\mathbf{w})$, where $\tilde{H}_i(\mathbf{w}) = \max_{\mathbf{y} \in \mathcal{Y}_i} L_i(\mathbf{y}) - \langle \mathbf{w}, \boldsymbol{\psi}_i(\mathbf{y}) \rangle$ is the structured hinge loss (using the notation from Lacoste-Julien et al. (2013)). We only assume access to the linear oracle computing $\tilde{H}_i(\mathbf{w})$. Let M_i have $(\boldsymbol{\psi}_i(\mathbf{y}))_{\mathbf{y} \in \mathcal{Y}_i}$ as columns. We can rewrite the minimization problem as a bilinear saddle point problem:

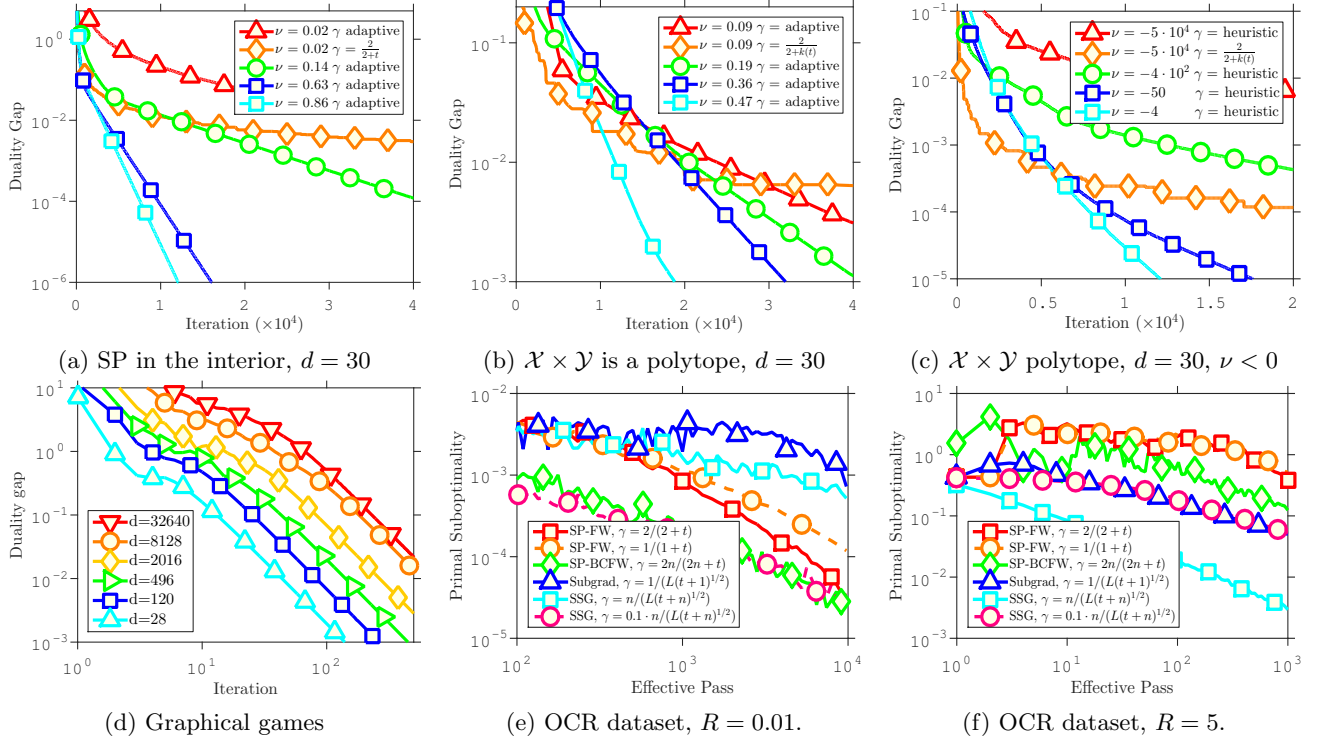


Figure 1: On Figure 1a, 1b and 1c, we plot on a semilog scale the best gap observed $\min_{i \leq t} g_i$ as a function of t . For experiments 1d, 1e and 1f, the objective function is bilinear and the convergence is sublinear. We give more details about these experiments in Appendix F.

$$\begin{aligned}
 & \min_{\|\mathbf{w}\|_1 \leq R} \frac{1}{n} \sum_i \left(\max_{\mathbf{y}_i \in \mathcal{Y}_i} \mathbf{L}_i^\top \mathbf{y}_i - \mathbf{w}^\top M_i \mathbf{y}_i \right) \\
 &= \min_{\|\mathbf{w}\|_1 \leq R} \frac{1}{n} \sum_i \left(\max_{\boldsymbol{\alpha}_i \in \Delta(|\mathcal{Y}_i|)} \mathbf{L}_i^\top \boldsymbol{\alpha}_i - \mathbf{w}^\top M_i \boldsymbol{\alpha}_i \right).
 \end{aligned} \quad (31)$$

Projecting onto $\Delta(|\mathcal{Y}_i|)$ is normally intractable as the size of $|\mathcal{Y}_i|$ is exponential, but the linear oracle is tractable by assumption. We performed experiments with 100 examples from the OCR dataset ($d_\omega = 4028$) (Taskar et al., 2003). We encoded the structure \mathcal{Y}_i of the i^{th} word with a Markov model: its k^{th} character $\mathcal{Y}_i^{(k)}$ only depends on \mathcal{Y}_i^{k-1} and \mathcal{Y}_i^{k+1} . In this case, the oracle function is simply the Viterbi algorithm Viterbi (1967). The average length of a word is approximately 8, hence the dimension of \mathcal{Y}_i is $d_{\mathcal{Y}_i} \approx 26^2 \cdot 8 = 5408$ leading to a large dimension for \mathcal{Y} , $d_{\mathcal{Y}} := \sum_{i=1}^n d_{\mathcal{Y}_i} \approx 5 \cdot 10^5$. We run the SP-FW algorithm with step size $\gamma_t = 1/(1+t)$ for which we have a convergence proof (Corollary 5), and with $\gamma_t = 2/(2+t)$, which normally gives better results for FW optimization. We compare with the projected subgradient method (projecting on the ℓ_1 -ball is tractable here) with step size $O(1/\sqrt{t})$ (the subgradient of $\tilde{H}_i(\mathbf{w})$ is $-\psi_i(\mathbf{y}_i^*)$). Following Lacoste-Julien et al. (2013), we also implement a block-coordinate (SP-BCFW) version of SP-FW and compare it with the stochastic projected subgradient method (SSG). As some of the algorithms

only work on the primal and to make our result comparable to Lacoste-Julien et al. (2013), we choose to plot the primal suboptimality error $p(\mathbf{w}_t) - p^*$ for the different algorithms in Figure 1e and 1f (the $\boldsymbol{\alpha}_t$ iterates for the SP approaches are thus ignored in this error). The performance of SP-BCFW is similar to SSG when we regularize the learning problem heavily (Figure 1e). However, under lower regularization (Figure 1f), SSG (with the correct step size scaling) is faster. This is consistent with the fact that $\boldsymbol{\alpha}_t \neq \boldsymbol{\alpha}^*$ implies larger errors on the primal suboptimality for the SP methods, but we note that an advantage of the SP-FW approach is that the scale of the step size is automatically chosen.

Conclusion. We proposed FW-style algorithms for saddle-point optimization with the same attractive properties as FW, in particular only requiring access to a LMO. We gave the first convergence result for a FW-style algorithm towards a saddle point over polytopes by building on the recent developments on the linear convergence analysis of AFW. However, our experiments let us believe that the condition $\nu > 0$ is not required for the convergence of FW-style algorithms. We thus conjecture that a refined convergence analysis could yield a linear rate for the general uniformly strongly convex-concave functions in both cases (I) and (P), paving the way for further theoretical work.

Acknowledgments

Thanks to N. Ruoizzi and A. Benchaouine for helpful discussions. Work supported in part by DARPA N66001-15-2-4026, N66001-15-C-4032 and NSF III-1526914, IIS-1451500, CCF-1302269.

References

- A. Ahmadinejad, S. Dehghani, Hajiaghayi, B. Lucier, H. Mahini, and S. Seddighin. From duels to battlefields: Computing equilibria of Blotto and other games. In *AAAI*, 2016.
- F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. Optimization with sparsity-inducing penalties. *Foundations and Trends in Machine Learning*, 2012.
- F. Barahona. On the computational complexity of Ising spin glass models. *J. Phys. A: Math. Gen.*, 1982.
- I. Bárány, S. Vempala, and A. Vetta. Nash equilibria in random games. *Random Structures & Algorithms*, 31(4):391–405, 2007.
- G. Brown. Iterative solution of games by fictitious play. *Activity analysis of production & allocation*, 1951.
- B. Cox, A. Juditsky, and A. Nemirovski. Decomposition techniques for bilinear saddle point problems and variational inequalities with affine monotone operators on domains given by linear minimization oracles. *arXiv preprint arXiv:1506.02444*, 2015.
- C. Daskalakis and Q. Pan. A counter-example to Karlin’s strong conjecture for fictitious play. In *FOCS*, 2014.
- V. F. Demyanov and A. M. Rubinov. *Approximate methods in optimization problems*. Elsevier, 1970.
- J. C. Dunn. Rates of convergence for conditional gradient algorithms near singular and nonsingular extremals. *SIAM Journal on Control and Optimization*, 1979.
- J. Edmonds. Paths, trees and flowers. *Canadian Journal of Mathematics*, 1965.
- M. Frank and P. Wolfe. An algorithm for quadratic programming. *Nav. Res. Logist. Q.*, 1956.
- D. Garber and E. Hazan. Faster rates for the Frank-Wolfe method over strongly-convex sets. In *ICML*, 2015.
- J. H. Hammond. *Solving asymmetric variational inequality problems and systems of equations with generalized nonlinear programming algorithms*. PhD thesis, Massachusetts Institute of Technology, 1984.
- J.-B. Hiriart-Urruty and C. Lemaréchal. *Convex analysis and minimization algorithms I: Fundamentals*. Springer, 2013.
- M. Jaggi. *Sparse convex optimization methods for machine learning*. PhD thesis, ETH Zürich, 2011.
- M. Jaggi. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In *ICML*, 2013.
- A. Juditsky and A. Nemirovski. Solving variational inequalities with monotone operators on domains given by linear minimization oracles. *Mathematical Programming*, 2016.
- S. Karlin. *Mathematical methods and theory in games, programming and economics*, 1960.
- J. Kelner, Y. Lee, L. Orrechia, and A. Sidford. An almost-linear-time algorithm for approximate max flow in undirected graphs, and its multicommodity generalizations. In *SODA*, 2014.
- S.-J. Kim, A. Magnani, and S. Boyd. Robust Fisher discriminant analysis. In *NIPS*, 2005.
- D. Koller, N. Megiddo, and B. Von Stengel. Fast algorithms for finding randomized strategies in game trees. In *STOC*, 1994.
- G. Korpelevich. The extragradient method for finding saddle points and other problems. *Matecon*, 1976.
- S. Lacoste-Julien and M. Jaggi. On the global linear convergence of Frank-Wolfe optimization variants. In *NIPS*, 2015.
- S. Lacoste-Julien, M. Jaggi, M. Schmidt, and P. Pletscher. Block-coordinate Frank-Wolfe optimization for structural SVMs. In *ICML*, 2013.
- G. Lan. The complexity of large-scale convex programming under a linear optimization oracle. *arXiv preprint arXiv:1309.5550*, 2013.
- T. Larsson and M. Patriksson. A class of gap functions for variational inequalities. *Math. Prog.*, 1994.
- E. S. Levitin and B. T. Polyak. Constrained minimization methods. *USSR Computational mathematics and mathematical physics*, 1966.
- A. F. Martins, N. A. Smith, P. M. Aguiar, and M. A. Figueiredo. Structured sparsity in structured prediction. In *EMNLP*, 2011.
- Y. Nesterov. Dual extrapolation and its applications to solving variational inequalities and related problems. *Mathematical Programming*, 2007.
- M. Patriksson. *Nonlinear Programming and Variational Inequality Problems: A Unified Approach*. Springer, 1999.
- B. T. Polyak. Existence theorems and convergence of minimizing sequences in extremum problems with restrictions. *Soviet Math. Dokl*, 1966.
- J. Robinson. An iterative method of solving a game. *Annals of mathematics*, 1951.

- H. N. Shapiro. Note on a computation method in the theory of games. *Communications on Pure and Applied Mathematics*, 1958.
- B. Taskar, C. Guestrin, and D. Koller. Max-margin markov networks. In *NIPS*, 2003.
- B. Taskar, S. Lacoste-Julien, and M. I. Jordan. Structured prediction, dual extragradient and Bregman projections. *Journal of Machine Learning Research*, 2006.
- J.-P. Vial. Strong and weak convexity of sets and functions. *Mathematics of Operations Research*, 1983.
- A. Viterbi. Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. *IEEE transactions on Information Theory*, 13(2):260–269, 1967.
- J. Von Neumann and O. Morgenstern. *Theory of games and economic behavior*. Princeton press, 1944.
- N. Xiu and J. Zhang. Some recent advances in projection-type methods for variational inequalities. *Journal of Computational and Applied Mathematics*, 2003.
- D. L. Zhu and P. Marcotte. Convergence properties of feasible descent methods for solving variational inequalities in banach spaces. *Computational Optimization and Applications*, 1998.

Appendix

Outline. Appendix A provides more details about the saddle point away-step Frank-Wolfe (SP-AFW) algorithm. Appendix B is about the affine invariant formulation of our algorithms, therein, we introduce some affine invariant constants and prove relevant bounds. Appendix C presents some relationships between the primal suboptimality and dual gaps useful for the convergence proof. Appendix D gives the affine invariant convergence proofs of SP-FW and SP-AFW in the strongly convex function setting introduced in Section 3. Appendix E gives the proof of linear convergence of SP-FW in the strongly convex set setting as defined in Section 4. Finally, Appendix F provides details on the experiments.

A Saddle point away-step Frank-Wolfe (SP-AFW)

In this section, we describe our algorithms SP-AFW and SP-PFW with a main focus on how the away direction is chosen. We also rigorously define a *drop step* and prove an upper bound on their number. In this section, we will assume that there exist two finite sets \mathcal{A} and \mathcal{B} such that $\mathcal{X} = \text{conv}(\mathcal{A})$ and $\mathcal{Y} = \text{conv}(\mathcal{B})$.

Active sets. Our definition of *active set* is an extension of the one provided in Lacoste-Julien and Jaggi (2015), we follow closely their notation and their results. Assume that we have the current expansion,

$$\mathbf{x}^{(t)} = \sum_{\mathbf{v}_x \in \mathcal{S}_x^{(t)}} \alpha_{\mathbf{v}_x}^{(t)} \mathbf{v}_x \quad \text{where} \quad \mathcal{S}_x^{(t)} := \left\{ \mathbf{v}_x \in \mathcal{A} ; \alpha_{\mathbf{v}_x}^{(t)} > 0 \right\}, \quad (32)$$

and a similar one for $\mathbf{y}^{(t)}$. Then, the current iterate has a sparse representation as a convex combination of all possible pairs of atoms belonging to $\mathcal{S}_x^{(t)}$ and $\mathcal{S}_y^{(t)}$, i.e.

$$\mathbf{z}^{(t)} = \sum_{\mathbf{v} \in \mathcal{S}^{(t)}} \alpha_{\mathbf{v}}^{(t)} \mathbf{v} \quad \text{where} \quad \mathcal{S}^{(t)} := \left\{ \mathbf{v} \in \mathcal{A} \times \mathcal{B} ; \alpha_{\mathbf{v}}^{(t)} := \alpha_{\mathbf{v}_x}^{(t)} \alpha_{\mathbf{v}_y}^{(t)} > 0 \right\}. \quad (33)$$

The set $\mathcal{S}^{(t)}$ is the current (implicit) *active set*. Note that after t iteration, the current iterate $\mathbf{z}^{(t)}$ is t -sparse whereas the size of the active set $\mathcal{S}^{(t)}$ defined in (33) can be t^2 . Fortunately, we only need to track at most t corners in \mathcal{A} and t ones in \mathcal{B} to get this bigger active set. We can now define the maximal step size for an away direction.

Maximal step size. For the standard AFW algorithm, Lacoste-Julien and Jaggi (2015) suggest to use the maximum step size $\gamma_{\max} = \alpha_{\mathbf{v}^{(t)}} / (1 - \alpha_{\mathbf{v}^{(t)}})$ when using the away direction $\mathbf{z}^{(t)} - \mathbf{v}^{(t)}$, to guarantee that the next iterate stays feasible. Because we have a product structure of two blocks, we actually consider more possible away directions by maintaining a separate convex combination on each block in our Algorithm 3 (SP-AFW) and 4 (SP-PFW). More precisely, suppose that we have $\mathbf{x}^{(t)} = \sum_{\mathbf{v}_x \in \mathcal{S}_x^{(t)}} \alpha_{\mathbf{v}_x}^{(t)} \mathbf{v}_x$ and $\mathbf{y}^{(t)} = \sum_{\mathbf{v}_y \in \mathcal{S}_y^{(t)}} \alpha_{\mathbf{v}_y}^{(t)} \mathbf{v}_y$, then the following maximum step size γ_{\max} (for AFW) ensures that the iterate $\mathbf{z}^{(t+1)}$ stays feasible:

$$\mathbf{z}^{(t+1)} := \mathbf{z}^{(t)} + \gamma_t \mathbf{d}_A^{(t)} \quad \text{with} \quad \gamma_t \in [0, \gamma_{\max}] \quad \text{and} \quad \gamma_{\max} := \min \left\{ \frac{\alpha_{\mathbf{v}_x}^{(t)}}{1 - \alpha_{\mathbf{v}_x}^{(t)}}, \frac{\alpha_{\mathbf{v}_y}^{(t)}}{1 - \alpha_{\mathbf{v}_y}^{(t)}} \right\}. \quad (34)$$

A larger γ_t makes one of the coefficients in the convex combination for the iterate negative, thus no more guaranteeing that the iterate stays feasible. A similar argument can be used to derive the maximal step size for the PFW direction in Algorithm 4.

Drop steps. A *drop step* is when $\gamma_t = \gamma_{\max}$ for the away-step update (34) (Lacoste-Julien and Jaggi, 2015). In this case, at least one corner is removed from the active set. We show later in Lemma 23 that we can still guarantee progress for this step, i.e. $w_{t+1} < w_t$, but this progress be arbitrarily small since γ_{\max} can be arbitrarily small. Lacoste-Julien and Jaggi (2015) shows that the number of drop steps for AFW is at most half of the number of iterations. Because we are maintaining two independent active sets in our formulation, we can obtain more drop steps, but we can still adapt their argument to obtain that the number of drop steps for SP-AFW is at most two thirds the number of iterations (assuming that the algorithm is initialized with only one atom per active set). In the SP-AFW algorithm, either a FW step is jointly made on both blocks, or an away-step is done on both blocks. Let us call A_t the number of FW steps (which potentially adds an atom in $S_x^{(t)}$ and $S_y^{(t)}$) and $D_t^{(x)}$ (resp $D_t^{(y)}$) the number of steps that removed at least one atom from $S_x^{(t)}$ ($S_y^{(t)}$). Finally, we call D_t the number of *drop steps*, i.e., the number of *away steps* where at least one atom from $S_x^{(t)}$ or $S_y^{(t)}$ have been removed (and thus $\gamma_t = \gamma_{\max}$ for these). Because a step is either a FW step or an away step, we have:

$$A_t + D_t \leq t. \quad (35)$$

We also have that $D_t^{(x)} + D_t^{(y)} \geq D_t$ by definition of D_t . Because a FW step adds at most one atom in an active set while a drop step removes one, we have (supposing that $|S_x^{(0)}| = |S_y^{(0)}| = 1$):

$$1 + A_t - D_t^{(x)} \geq |S_x^{(t)}| \quad \text{and} \quad 1 + A_t - D_t^{(y)} \geq |S_y^{(t)}|. \quad (36)$$

Adding these two relations, we get:

$$2 + 2A_t \geq |S_x^{(t)}| + |S_y^{(t)}| + D_t^{(x)} + D_t^{(y)} \geq 2 + D_t, \quad (37)$$

using the fact that each active set as at least one element. We thus obtain $D_t \leq 2A_t$. Combining with (35), we get:

$$D_t \leq \frac{2}{3}t, \quad (38)$$

as claimed.

B Affine invariant formulation of SP-FW

In this section, we define the affine invariant constants of a convex function f and their extension to a convex-concave function \mathcal{L} . These constants are important as the FW-type algorithms are affine invariant if their step size are defined using affine invariant quantities. We can upper bound these constants using the non affine invariant constants defined in the main paper. Hence a convergence rate with affine invariant constants will immediately imply a rate with the constant introduced in the main paper.

B.1 The Lipschitz constants

We define the Lipschitz constant L of the gradient of the function f with respect to the norm $\|\cdot\|$ by using a dual pairing of norms, i.e. L is a constant such that

$$\forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}, \quad \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}')\|_* \leq L\|\mathbf{x} - \mathbf{x}'\|, \quad (39)$$

where $\|\mathbf{y}\|_* := \sup_{\mathbf{x} \in \mathbb{R}^d, \|\mathbf{x}\| \leq 1} \mathbf{y}^T \mathbf{x}$ is the dual norm of $\|\cdot\|$. For a convex-concave function, we also consider the partial Lipschitz constants with respect to different blocks as follows.

For more generality, we consider the dual pairing of norms $(\|\cdot\|_{\mathcal{X}}, \|\cdot\|_{\mathcal{X}^*})$ on \mathcal{X} , and similarly $(\|\cdot\|_{\mathcal{Y}}, \|\cdot\|_{\mathcal{Y}^*})$ on \mathcal{Y} . We also define the norm on the product space $\mathcal{X} \times \mathcal{Y}$ as the ℓ_1 -norm on the components: $\|(\mathbf{x}, \mathbf{y})\|_{\mathcal{X} \times \mathcal{Y}} := \|\mathbf{x}\|_{\mathcal{X}} + \|\mathbf{y}\|_{\mathcal{Y}}$. We thus have that the dual norm of $\mathcal{X} \times \mathcal{Y}$ is the

ℓ_∞ -norm of the dual norms: $\|(\mathbf{x}, \mathbf{y})\|_{(\mathcal{X} \times \mathcal{Y})^*} = \max(\|\mathbf{x}\|_{\mathcal{X}^*}, \|\mathbf{y}\|_{\mathcal{Y}^*})$. The *partial* Lipschitz constants L_{XX}, L_{YY}, L_{XY} and L_{YX} of the gradient of the function \mathcal{L} with respect to these norms are the constants such that for all $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ and $\mathbf{y}, \mathbf{y}' \in \mathcal{Y}$,

$$\begin{aligned} \|\nabla_x \mathcal{L}(\mathbf{x}, \mathbf{y}) - \nabla_x \mathcal{L}(\mathbf{x}', \mathbf{y})\|_{\mathcal{X}^*} &\leq L_{XX} \|\mathbf{x} - \mathbf{x}'\|_{\mathcal{X}}, & \|\nabla_y \mathcal{L}(\mathbf{x}, \mathbf{y}) - \nabla_y \mathcal{L}(\mathbf{x}, \mathbf{y}')\|_{\mathcal{Y}^*} &\leq L_{YY} \|\mathbf{y} - \mathbf{y}'\|_{\mathcal{Y}}, \\ \|\nabla_x \mathcal{L}(\mathbf{x}, \mathbf{y}) - \nabla_x \mathcal{L}(\mathbf{x}, \mathbf{y}')\|_{\mathcal{X}^*} &\leq L_{XY} \|\mathbf{y} - \mathbf{y}'\|_{\mathcal{Y}}, & \|\nabla_y \mathcal{L}(\mathbf{x}, \mathbf{y}) - \nabla_y \mathcal{L}(\mathbf{x}', \mathbf{y})\|_{\mathcal{Y}^*} &\leq L_{YX} \|\mathbf{x} - \mathbf{x}'\|_{\mathcal{X}}. \end{aligned} \quad (40)$$

Note that the cross partial Lipschitz constants L_{XY} and L_{YX} do not necessarily use a dual pairing as \mathcal{X} and \mathcal{Y} could be very different spaces. On the other hand, as the possibilities in (40) are special cases of (39) when considering the ℓ_1 -norm of this product domain, one can easily deduce that the partial Lipschitz constants can always be taken to be smaller than the full Lipschitz constant for the gradient of \mathcal{L} , i.e., we have that $L \geq \max(L_{XX}, L_{XY}, L_{YX}, L_{YY})$.

B.2 The curvature: an affine invariant measure of smoothness

To prove the convergence of the Frank-Wolfe algorithm, the typical affine invariant analysis proof in the FW literature assumes that the curvature of the objective function is bounded, where the curvature is defined by Jaggi (2013) for example. We give below a slight generalization of this curvature notion in order to handle the convergence analysis of FW with away-steps.² It has the same upper bound as the traditional curvature constant (see Proposition 6).

Curvature. [Slight generalization of Jaggi (2013)] Let $f : \mathcal{X} \rightarrow \mathbb{R}$ be a convex function, we define the curvature C_f of f as

$$C_f := \sup_{\substack{\mathbf{x}, \mathbf{s}, \mathbf{v} \in \mathcal{X}, \\ \gamma > 0 \text{ s.t.} \\ \mathbf{x}_\gamma := \mathbf{x} + \gamma \mathbf{d} \in \mathcal{X} \\ \text{with } \mathbf{d} := \mathbf{s} - \mathbf{v}}} \frac{2}{\gamma^2} (f(\mathbf{x}_\gamma) - f(\mathbf{x}) - \gamma \langle \mathbf{d}, \nabla f(\mathbf{x}) \rangle). \quad (41)$$

Note that only the *feasible* step sizes γ are considered in the definition of C_f , i.e., γ such that $\mathbf{x}_\gamma \in \mathcal{X}$. If the gradient of the objective function is Lipschitz continuous, the curvature is upper bounded.

Proposition 6 (Simple generalization of Lemma 7 in Jaggi (2013)). *Let f be a convex and continuously differentiable function on \mathcal{X} with its gradient ∇f L -Lipschitz continuous w.r.t. some norm $\|\cdot\|$ in dual pairing over the domain \mathcal{X} . Then*

$$C_f \leq D_{\mathcal{X}}^2 L, \quad (42)$$

where $D_{\mathcal{X}} := \sup_{\mathbf{x}, \mathbf{x}' \in \mathcal{X}} \|\mathbf{x} - \mathbf{x}'\|$ is the diameter of \mathcal{X} .

Lemma 1.2.3 in Nesterov (2004), Jaggi (2013). Let $\mathbf{x}, \mathbf{s}, \mathbf{v} \in \mathcal{X}$, set $\mathbf{d} := \mathbf{s} - \mathbf{v}$ and $\mathbf{x}_\gamma = \mathbf{x} + \gamma \mathbf{d}$ for some $\gamma > 0$ such that $\mathbf{x}_\gamma \in \mathcal{X}$. Then by the fundamental theorem of calculus,

$$f(\mathbf{x}_\gamma) = f(\mathbf{x}) + \int_0^\gamma \langle \mathbf{d}, \nabla f(\mathbf{x} + t\mathbf{d}) \rangle dt. \quad (43)$$

²The change is to consider the more general directions $\mathbf{s} - \mathbf{v}$ instead of just $\mathbf{s} - \mathbf{x}$, and also any feasible positive step size. See also Footnote 8 in Lacoste-Julien and Jaggi (2013) for a related discussion. A different (bigger) constant was required in (Lacoste-Julien and Jaggi, 2015) for the analysis of AFW because they used a line-search.

Hence, we can write

$$\begin{aligned}
 f(\mathbf{x}_\gamma) - f(\mathbf{x}) - \gamma \langle \mathbf{d}, \nabla f(\mathbf{x}) \rangle &= \int_0^\gamma \langle \mathbf{d}, \nabla f(\mathbf{x} + t\mathbf{d}) - \nabla f(\mathbf{x}) \rangle dt \\
 &\leq \|\mathbf{d}\| \int_0^\gamma \|\nabla f(\mathbf{x} + t\mathbf{d}) - \nabla f(\mathbf{x})\|_* dt \\
 &\leq D_{\mathcal{X}}^2 L \int_0^\gamma t dt \\
 &\leq \frac{\gamma^2}{2} D_{\mathcal{X}}^2 L.
 \end{aligned} \tag{44}$$

Thus for all $\mathbf{x}, \mathbf{s}, \mathbf{v} \in \mathcal{X}$ and $\mathbf{x}_\gamma = \mathbf{x} + \gamma(\mathbf{s} - \mathbf{v})$ for $\gamma > 0$ such that $\mathbf{x}_\gamma \in \mathcal{X}$, we have

$$\frac{2}{\gamma^2} (f(\mathbf{x}_\gamma) - f(\mathbf{x}) - \gamma \langle \mathbf{s} - \mathbf{v}, \nabla f(\mathbf{x}) \rangle) \leq LD_{\mathcal{X}}^2. \tag{45}$$

The supremum is then upper bounded by the claimed quantity. \square

Osokin et al. (2016, Appendix C.1) illustrate well the importance of the affine invariant curvature constant for Frank-Wolfe algorithms in their paragraph titled ‘‘Lipschitz and curvature constants’’. They provide a concrete example where the wrong choice of norm for a specific domain \mathcal{X} can make the upper bound of Proposition 6 extremely loose, and thus practically useless for an analysis.

We will therefore extend the curvature constant to the convex-concave function \mathcal{L} by simply defining it as the maximum of the curvatures of the functions belonging to the family $(\mathbf{x}' \mapsto \mathcal{L}(\mathbf{x}', \mathbf{y}), \mathbf{y}' \mapsto -\mathcal{L}(\mathbf{x}, \mathbf{y}'))_{\mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}}$ (see Section B.4). But before that, we review affine invariant analogues of the strong convexity constants that will be useful for the analysis.

B.3 Affine invariant measures of strong convexity

In this section, we review two affine invariant measures of strong convexity that were proposed by Lacoste-Julien and Jaggi (2013) (Lacoste-Julien and Jaggi, 2015) for the affine invariant linear convergence analysis of the standard Frank-Wolfe algorithm (using the ‘‘interior strong convexity constant’’) or the away-step Frank-Wolfe algorithm (using the ‘‘geometric strong convexity constant’’). We will re-use them for the affine invariant analysis of the convergence of SP-FW or SP-AFW algorithms. In a similar way as the curvature constant C_f includes information about the constraint set \mathcal{X} and the Lipschitz continuity of the gradient of f together, these constants both include the information about the constraint set \mathcal{X} and the strong convexity of a function f together.

Interior strong convexity constant. [based on Lacoste-Julien and Jaggi (2013)] Let \mathbf{x}_c be a point in the relative interior of \mathcal{X} . The *interior strong convexity constant* for f with respect to the reference point \mathbf{x}_c is defined as

$$\mu_f^{\mathbf{x}_c} := \inf_{\substack{\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{x}_c\} \\ \mathbf{s} = \bar{\mathbf{s}}(\mathbf{x}, \mathbf{x}_c, \mathcal{X}) \\ \gamma \in (0, 1], \\ \mathbf{z} = \mathbf{x} + \gamma(\mathbf{s} - \mathbf{x})}} \frac{2}{\gamma^2} (f(\mathbf{z}) - f(\mathbf{x}) - \langle \mathbf{z} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle). \tag{46}$$

Here, we follow the notation of Lacoste-Julien and Jaggi (2013) and take the point \mathbf{s} to be the point where the ray from \mathbf{x} to the reference point \mathbf{x}_c pinches the boundary of the set \mathcal{X} , i.e. $\bar{\mathbf{s}}(\mathbf{x}, \mathbf{x}_c, \mathcal{X}) := \text{ray}(\mathbf{x}, \mathbf{x}_c) \cap \partial\mathcal{X}$, where $\partial\mathcal{X}$ is the boundary of the convex set \mathcal{X} .

We note that in the original definition (Lacoste-Julien and Jaggi, 2013), \mathbf{x}_c was the (unique) optimum point for a strongly convex function f over \mathcal{X} . The optimality of \mathbf{x}_c is actually not needed in the

definition and so we generalize it here to any point \mathbf{x}_c in the relative interior of \mathcal{X} , as this will be useful in our convergence proof for SP-FW.

For completeness, we include here the important lower bound from (Lacoste-Julien and Jaggi, 2013) on the interior strong convexity constant in terms of the strong convexity of the function f .

Proposition 7 (Lower bound on $\mu^{\mathbf{x}_c}$ from Lacoste-Julien and Jaggi (2013, Lemma 2)). *Let f be a convex differentiable function and suppose that f is strongly convex w.r.t. to some arbitrary norm $\|\cdot\|$ over the domain \mathcal{X} with strong-convexity constant $\mu_f > 0$. Furthermore, suppose that the reference point \mathbf{x}_c lies in the relative interior of \mathcal{X} , i.e., $\delta_c := \min_{\mathbf{s} \in \partial\mathcal{X}} \|\mathbf{s} - \mathbf{x}_c\| > 0$. Then the interior strong convexity constant $\mu_f^{\mathbf{x}_c}$ (46) is lower bounded as follows:*

$$\mu_f^{\mathbf{x}_c} \geq \mu_f \delta_c^2. \quad (47)$$

Proof. Let \mathbf{x} and \mathbf{z} be defined as in (46), i.e., $\mathbf{z} = \mathbf{x} + \gamma(\mathbf{s} - \mathbf{x})$ for some $\gamma > 0$ and where \mathbf{s} intersects the boundary of \mathcal{X} with the ray going from \mathbf{x} to \mathbf{x}_c . By the strong convexity of f , we have

$$f(\mathbf{z}) - f(\mathbf{x}) - \langle \mathbf{z} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle \geq \|\mathbf{z} - \mathbf{x}\|^2 \frac{\mu_f}{2} = \gamma^2 \|\mathbf{s} - \mathbf{x}\|^2 \frac{\mu_f}{2}. \quad (48)$$

From the definition of \mathbf{s} , we have that \mathbf{x}_c lies between \mathbf{x} and \mathbf{s} and thus: $\|\mathbf{s} - \mathbf{x}\| \geq \|\mathbf{s} - \mathbf{x}_c\| \geq \delta_c$. Combining with (48), we conclude

$$f(\mathbf{z}) - f(\mathbf{x}) - \langle \mathbf{z} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle \geq \gamma^2 \delta_c^2 \frac{\mu_f}{2}, \quad (49)$$

and therefore

$$\mu_f^{\mathbf{x}_c} \geq \delta_c^2 \mu_f. \quad (50)$$

□

We now present the affine invariant constant used in the global linear convergence analysis of Frank-Wolfe variants when the convex set \mathcal{X} is a polytope. The *geometric strong convexity constant* was originally introduced by Lacoste-Julien and Jaggi (2013) and (Lacoste-Julien and Jaggi, 2015). To avoid any ambiguity, we will re-use their definitions verbatim in the rest of this section, starting first with a few geometrical definitions and then presenting the affine invariant constant. In these definitions, they assume that a *finite* set \mathcal{A} of vectors (that they call *atoms*) is given such that $\mathcal{X} = \text{conv}(\mathcal{A})$ (which always exists when \mathcal{X} is a polytope).

Directional Width. [Lacoste-Julien and Jaggi (2015)] The *directional width* of a set \mathcal{A} with respect to a direction \mathbf{r} is defined as $\text{dir}W(\mathcal{A}, \mathbf{r}) := \max_{\mathbf{s}, \mathbf{v} \in \mathcal{A}} \left\langle \frac{\mathbf{r}}{\|\mathbf{r}\|_2}, \mathbf{s} - \mathbf{v} \right\rangle$. The *width* of \mathcal{A} is the minimum directional width over all possible directions in its affine hull.

Pyramidal Directional Width. [Lacoste-Julien and Jaggi (2015)] We define the *pyramidal directional width* of a set \mathcal{A} with respect to a direction \mathbf{r} and a base point $\mathbf{x} \in \mathcal{X}$ to be

$$P\text{dir}W(\mathcal{A}, \mathbf{r}, \mathbf{x}) := \min_{\mathcal{S} \in \mathcal{S}_{\mathbf{x}}} \text{dir}W(\mathcal{S} \cup \{\mathbf{s}(\mathcal{A}, \mathbf{r})\}, \mathbf{r}) = \min_{\mathcal{S} \in \mathcal{S}_{\mathbf{x}}} \max_{\mathbf{s} \in \mathcal{A}, \mathbf{v} \in \mathcal{S}} \left\langle \frac{\mathbf{r}}{\|\mathbf{r}\|_2}, \mathbf{s} - \mathbf{v} \right\rangle, \quad (51)$$

where $\mathcal{S}_{\mathbf{x}} := \{\mathcal{S} \mid \mathcal{S} \subseteq \mathcal{A} \text{ such that } \mathbf{x} \text{ is a proper}^3 \text{ convex combination of all the elements in } \mathcal{S}\}$, and $\mathbf{s}(\mathcal{A}, \mathbf{r}) := \arg \max_{\mathbf{v} \in \mathcal{A}} \langle \mathbf{r}, \mathbf{v} \rangle$ is the FW atom used as a summit, when using the convention in this section that $\mathbf{r} := -\nabla f(\mathbf{x})$.

³By *proper* convex combination, we mean that all coefficients are non-zero in the convex combination.

Pyramidal Width. [Lacoste-Julien and Jaggi (2015)] To define the pyramidal width of a set, we take the minimum over the cone of possible *feasible* directions \mathbf{r} (in order to avoid the problem of zero width).

A direction \mathbf{r} is *feasible* for \mathcal{A} from \mathbf{x} if it points inwards $\text{conv}(\mathcal{A})$, (i.e. $\mathbf{r} \in \text{cone}(\mathcal{A} - \mathbf{x})$).

We define the *pyramidal width* of a set \mathcal{A} to be the smallest pyramidal width of all its faces, i.e.

$$PWidth(\mathcal{A}) := \min_{\substack{\mathcal{K} \in \text{faces}(\text{conv}(\mathcal{A})) \\ \mathbf{x} \in \mathcal{K} \\ \mathbf{r} \in \text{cone}(\mathcal{K} - \mathbf{x}) \setminus \{\mathbf{0}\}}} PdirW(\mathcal{K} \cap \mathcal{A}, \mathbf{r}, \mathbf{x}). \quad (52)$$

Geometric strong convexity constant. [Lacoste-Julien and Jaggi (2015)] The *geometric strong convexity constant* of f (over the set of atoms \mathcal{A} which is left implicit) is:

$$\mu_f^{\mathcal{A}} := \inf_{\mathbf{x} \in \mathcal{X}} \inf_{\substack{\mathbf{x}^* \in \mathcal{X} \\ s.t. \langle \nabla f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle < 0}} \frac{2}{\gamma^{\mathcal{A}}(\mathbf{x}, \mathbf{x}^*)^2} (f(\mathbf{x}^*) - f(\mathbf{x}) - \langle \mathbf{x}^* - \mathbf{x}, \nabla f(\mathbf{x}) \rangle) \quad (53)$$

where $\gamma^{\mathcal{A}}(\mathbf{x}, \mathbf{x}^*) := \frac{\langle -\nabla f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle}{\langle -\nabla f(\mathbf{x}), \mathbf{s}_f(\mathbf{x}) - \mathbf{v}_f(\mathbf{x}) \rangle}$ and $\mathcal{X} = \text{conv}(\mathcal{A})$. The quantity $\mathbf{s}_f(\mathbf{x})$ represents the FW corner picked when running the FW algorithm on f when at \mathbf{x} ; while $\mathbf{v}_f(\mathbf{x})$ represents the worst-case possible away atom that AFW could pick (and this is where the dependence on \mathcal{A} appears). We now define these quantities more precisely. Recall that the set of possible active sets is $\mathcal{S}_{\mathbf{x}} := \{\mathcal{S} \mid \mathcal{S} \subseteq \mathcal{A} \text{ such that } \mathbf{x} \text{ is a proper convex combination of all the elements in } \mathcal{S}\}$. For a given set \mathcal{S} , we write $\mathbf{v}_{\mathcal{S}}(\mathbf{x}) := \arg \max_{\mathbf{v} \in \mathcal{S}} \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle$ for the away atom in the algorithm supposing that the current set of active atoms is \mathcal{S} . Finally, we define $\mathbf{v}_f(\mathbf{x}) := \arg \min_{\{\mathbf{v} = \mathbf{v}_{\mathcal{S}}(\mathbf{x}) \mid \mathcal{S} \in \mathcal{S}_{\mathbf{x}}\}} \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle$ to be the worst-case away

atom (that is, the atom which would yield the smallest away descent). An important property coming from this definition that we will use later is that for $\mathbf{s}^{(t)}$ and $\mathbf{v}^{(t)}$ being possible FW and away atoms (respectively) appearing during the AFW algorithm (consider Algorithm 3 ran only on \mathcal{X}), then we have:

$$g_t^{\text{PFW}} := \left\langle \mathbf{s}^{(t)} - \mathbf{v}^{(t)}, -\nabla f(\mathbf{x}^{(t)}) \right\rangle \geq \left\langle \mathbf{s}_f(\mathbf{x}^{(t)}) - \mathbf{v}_f(\mathbf{x}^{(t)}), -\nabla f(\mathbf{x}^{(t)}) \right\rangle. \quad (54)$$

The following important theorem from (Lacoste-Julien and Jaggi, 2015) lower bounds the geometric strong convexity constant of f in terms of both the strong convexity constant of f , as well as the pyramidal width of $\mathcal{X} = \text{conv}(\mathcal{A})$ defined as $PWidth(\mathcal{A})$ (52).

Proposition 8 (Lower bound for $\mu_f^{\mathcal{A}}$ from Lacoste-Julien and Jaggi (2015, Theorem 6)). *Let f be a convex differentiable function and suppose that f is μ -strongly convex w.r.t. to the Euclidean norm $\|\cdot\|_2$ over the domain $\mathcal{X} = \text{conv}(\mathcal{A})$ with strong-convexity constant $\mu \geq 0$. Then*

$$\mu_f^{\mathcal{A}} \geq \mu \cdot (PWidth(\mathcal{A}))^2. \quad (55)$$

The pyramidal width (52) is a geometric quantity with a somewhat intricate definition. Its value is still unknown for many sets (though always strictly positive for finite sets), but Lacoste-Julien and Jaggi (2015, Lemma 4) give its value for the unit cube in \mathbb{R}^d as $1/\sqrt{d}$.

B.4 Curvature and interior strong convexity constant for a convex-concave function

In this subsection, we propose simple convex-concave extensions of the definitions of the affine invariant constants defined introduced in the two previous sections.

To define the convex-concave curvature, we introduce the sets \mathcal{F} and \mathcal{G} of the marginal convex functions.

$$\mathcal{F} := \{\mathbf{x}' \mapsto \mathcal{L}(\mathbf{x}', \mathbf{y})\}_{\mathbf{y} \in \mathcal{Y}} \quad \text{and} \quad \mathcal{G} := \{\mathbf{y}' \mapsto -\mathcal{L}(\mathbf{x}, \mathbf{y}')\}_{\mathbf{x} \in \mathcal{X}}. \quad (56)$$

Let $\mathcal{L} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ a convex-concave function, we define the curvature pair $(C_{\mathcal{L}_x}, C_{\mathcal{L}_y})$ of \mathcal{L} as

$$(C_{\mathcal{L}_x}, C_{\mathcal{L}_y}) := \left(\sup_{f \in \mathcal{F}} C_f, \sup_{g \in \mathcal{G}} C_g \right). \quad (57)$$

and the curvature of \mathcal{L} as

$$C_{\mathcal{L}} := \frac{C_{\mathcal{L}_x} + C_{\mathcal{L}_y}}{2}. \quad (58)$$

An upper bound on this quantity follows directly from the upper bound on the convex case (Lemma 7 of Jaggi (2013), repeated in our Proposition 6) :

Proposition 9. *Let $\mathcal{L} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ be a differentiable convex-concave function. If \mathcal{X} and \mathcal{Y} are compact and $\nabla \mathcal{L}$ is Lipschitz continuous, then the curvature of \mathcal{L} is bounded by $\frac{1}{2}(L_{XX}D_{\mathcal{X}}^2 + L_{YY}D_{\mathcal{Y}}^2)$, where L_{XX} (resp L_{YY}) is the largest Lipschitz constant respect to \mathbf{x} (\mathbf{y}) of $\mathbf{x} \mapsto \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{y})$ ($\mathbf{y} \mapsto \nabla_{\mathbf{y}} \mathcal{L}(\mathbf{x}, \mathbf{y})$).*

Proof. Let f in \mathcal{F} ,

$$C_f \leq \text{Lip}(\nabla f) D_{\mathcal{X}}^2 \leq L_{XX} D_{\mathcal{X}}^2. \quad (59)$$

Similarly, let g in \mathcal{G} ,

$$C_g \leq \text{Lip}(\nabla g) D_{\mathcal{Y}}^2 \leq L_{YY} D_{\mathcal{Y}}^2. \quad (60)$$

Consequently,

$$C_{\mathcal{L}} = \frac{1}{2}(\sup_{f \in \mathcal{F}} C_f + \sup_{g \in \mathcal{G}} C_g) \leq \frac{1}{2}(L_{XX} D_{\mathcal{X}}^2 + L_{YY} D_{\mathcal{Y}}^2). \quad (61)$$

Where $D_{\mathcal{X}}$ and $D_{\mathcal{Y}}$ are the respective diameter of \mathcal{X} and \mathcal{Y} . □

Note that L_{XX} and L_{YY} are upper bounded by the global Lipschitz constant of $\nabla \mathcal{L}$. Similarly, we define various notions of strong convex-concavity in the following.

Uniform strong convex-concavity constant. The uniform strong convex-concavity constants is defined as

$$(\mu_{\mathcal{X}}, \mu_{\mathcal{Y}}) := \left(\inf_{f \in \mathcal{F}} \mu_f, \inf_{g \in \mathcal{G}} \mu_g \right) \quad (62)$$

where μ_f is the strong convexity constant of f and μ_g the strong convexity of g .

Under some assumptions this quantity is positive.

Proposition 10. *If the second derivative of \mathcal{L} is continuous, \mathcal{X} and \mathcal{Y} are compact and if for all $f \in \mathcal{F} \cup \mathcal{G}$, $\mu_f > 0$, then $\mu_{\mathcal{X}}$ and $\mu_{\mathcal{Y}}$ are positive.*

Proof. Let us introduce $H_x(\mathbf{x}, \mathbf{y}) := \nabla_{\mathbf{x}}^2 \mathcal{L}(\mathbf{x}, \mathbf{y})$ the Hessian of the function $\mathbf{x} \mapsto \mathcal{L}(\mathbf{x}, \mathbf{y})$. We want to show that the smallest eigenvalue is uniformly bounded on $\mathcal{X} \times \mathcal{Y}$. We know that the smallest eigenvalue lower bounds $\mu_{\mathcal{X}}$,

$$\mu_{\mathcal{X}} \geq \inf_{\substack{(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y} \\ \|\mathbf{u}\|_2 = 1}} \langle \mathbf{u}, H_x(\mathbf{x}, \mathbf{y}) \cdot \mathbf{u} \rangle. \quad (63)$$

But $H_x(\cdot)$ is continuous (because $\nabla_{\mathbf{x}}^2 \mathcal{L}(\cdot)$ is continuous by assumption) and then the function $(\mathbf{u}, \mathbf{x}, \mathbf{y}) \mapsto \langle \mathbf{u}, H_x(\mathbf{x}, \mathbf{y}) \cdot \mathbf{u} \rangle$ is continuous. Hence since $\mathcal{X} \times \mathcal{Y}$ and the unit ball are compact, the infimum is a minimum which can't be 0 by assumption. Hence $\mu_{\mathcal{X}}$ is positive. Doing the same thing with the smallest eigenvalue of $-\nabla_{\mathbf{y}}^2 \mathcal{L}(\mathbf{x}, \mathbf{y})$, we get that $\mu_{\mathcal{Y}} > 0$. □

A common family of saddle point objectives is of the form $f(x) + x^T My - g(y)$. In this case, we get simply that $(\mu_{\mathcal{X}}, \mu_{\mathcal{Y}}) = (\mu_f, \mu_g)$. An equivalent definition for the uniform strong convex-concavity constant is: \mathcal{L} is $(\mu_{\mathcal{X}}, \mu_{\mathcal{Y}})$ -uniform strongly convex-concave function if

$$(\mathbf{x}, \mathbf{y}) \mapsto \mathcal{L}(\mathbf{x}, \mathbf{y}) - \frac{\mu_{\mathcal{X}}}{2} \|\mathbf{x}\|^2 + \frac{\mu_{\mathcal{Y}}}{2} \|\mathbf{y}\|^2 \quad (64)$$

is convex-concave.

The following proposition relates the distance between the saddle point and the values of the function. It is a direct consequence from the uniform strong convex-concavity definition (62).

Proposition 11. *Let \mathcal{L} be a uniformly strongly convex-concave function and $(\mathbf{x}^*, \mathbf{y}^*)$ the saddle point of \mathcal{L} . Then we have for all \mathbf{x} in \mathcal{X} and $\mathbf{y} \in \mathcal{Y}$,*

$$\sqrt{\mathcal{L}(\mathbf{x}, \mathbf{y}^*) - \mathcal{L}^*} \geq \|\mathbf{x}^* - \mathbf{x}\| \sqrt{\frac{\mu_{\mathcal{X}}}{2}} \quad \text{and} \quad \sqrt{\mathcal{L}^* - \mathcal{L}(\mathbf{x}^*, \mathbf{y})} \geq \|\mathbf{y}^* - \mathbf{y}\| \sqrt{\frac{\mu_{\mathcal{Y}}}{2}}. \quad (65)$$

Proof. The saddle point $(\mathbf{x}^*, \mathbf{y}^*)$ is the optimal point of the two strongly convex functions $\mathbf{x} \mapsto \mathcal{L}(\mathbf{x}, \mathbf{y}^*)$ and the function $\mathbf{y} \mapsto -\mathcal{L}(\mathbf{x}^*, \mathbf{y})$, so we can use the property of strong convexity on each function and the fact that $\mu_{\mathcal{X}}$ lower bounds the strong convexity constant of $\mathcal{L}(\cdot, \mathbf{y}^*)$ (and similarly for $\mu_{\mathcal{Y}}$ with $-\mathcal{L}(\mathbf{x}^*, \cdot)$) as per the definition (62), to get the required conclusion. \square

Now we will introduce the uniform strong convex-concavity constants relatively to our saddle point.

Interior strong convex-concavity. The SP-FW interior strong convex-concavity constants (with respect to the reference point $(\mathbf{x}_c, \mathbf{y}_c)$) are defined as:

$$(\mu_{\mathcal{L}}^{\mathbf{x}_c}, \mu_{\mathcal{L}}^{\mathbf{y}_c}) := \left(\inf_{f \in \mathcal{F}} \mu_f^{\mathbf{x}_c}, \inf_{g \in \mathcal{G}} \mu_g^{\mathbf{y}_c} \right) \quad (66)$$

where $\mu_f^{\mathbf{x}_c}$ is the interior strong convexity constant of f w.r.t to the point \mathbf{x}_c and $\mu_g^{\mathbf{y}_c}$ is the interior strong convexity constant w.r.t to the point \mathbf{y}_c . The sets \mathcal{F} and \mathcal{G} are defined in (56). We also define the smallest quantity of both (with the reference point $(\mathbf{x}_c, \mathbf{y}_c)$ implicit):

$$\mu_{\mathcal{L}}^{\text{int}} = \min\{\mu_{\mathcal{L}}^{\mathbf{x}_c}, \mu_{\mathcal{L}}^{\mathbf{y}_c}\}. \quad (67)$$

We can lower bound this constant by a quantity depending on the uniform strong convexity constant and the distance of the saddle point to the boundary. The propositions on the strong convex-concavity directly follow from the previous definitions and the analogous proposition on the convex case (Proposition 7)

Proposition 12. *Let \mathcal{L} be a convex-concave function. If the reference point $(\mathbf{x}_c, \mathbf{y}_c)$ belongs to the relative interior of $\mathcal{X} \times \mathcal{Y}$ and if the function \mathcal{L} is strongly convex-concave with a strong convex-concavity constant $\mu > 0$, then $\mu_{\mathcal{L}}^{\text{int}}$ is lower bounded away from zero. More precisely, define $\delta_x := \min_{\mathbf{s}_x \in \partial \mathcal{X}} \|\mathbf{s}_x - \mathbf{x}_c\| > 0$ and $\delta_y := \min_{\mathbf{s}_y \in \partial \mathcal{Y}} \|\mathbf{s}_y - \mathbf{y}_c\|$. Then we have,*

$$\mu_{\mathcal{L}}^{\mathbf{x}_c} \geq \mu_{\mathcal{X}} \delta_x^2 \quad \text{and} \quad \mu_{\mathcal{L}}^{\mathbf{y}_c} \geq \mu_{\mathcal{Y}} \delta_y^2. \quad (68)$$

Proof. Using the Proposition 7 we have,

$$\mu_f^{\mathbf{x}_c} \geq \mu_f \cdot \delta_x^2 \geq \mu_{\mathcal{X}} \cdot \delta_x^2, \quad (69)$$

and,

$$\mu_g^{\mathbf{y}_c} \geq \mu_g \cdot \delta_y^2 \geq \mu_{\mathcal{Y}} \cdot \delta_y^2. \quad (70)$$

\square

When the saddle point is not in the interior of the domain, we define next a constant that takes in consideration the geometry of the sets. If the sets are polytopes, then this constant is positive.

Geometric strong convex-concavity. The SP-FW *geometric* strong convex-concavity constants are defined analogously as the interior strong convex-concavity constants,

$$\left(\mu_{\mathcal{L}_x}^{\mathbb{A}}, \mu_{\mathcal{L}_y}^{\mathbb{A}}\right) := \left(\min_{f \in \mathcal{F}} \mu_f^{\mathbb{A}}, \min_{g \in \mathcal{G}} \mu_g^{\mathbb{A}}\right); \quad \mu_{\mathcal{L}}^{\mathbb{A}} := \min\left(\mu_{\mathcal{L}_x}^{\mathbb{A}}, \mu_{\mathcal{L}_y}^{\mathbb{A}}\right), \quad (71)$$

where $\mu_f^{\mathbb{A}}$ is the geometric strong convexity constant of $f \in \mathcal{F}$ (over \mathcal{A}) as defined in (53) (and similarly $\mu_g^{\mathbb{A}}$ is the geometric strong convexity constant of $g \in \mathcal{G}$ over \mathcal{B}).

It is straightforward to notice that the lower bound on the geometric strong convexity constant (Proposition 8) can be extended to the geometric strong convex-concavity constants (where $\mu_{\mathcal{X}}$ and $\mu_{\mathcal{Y}}$ are now assumed to be defined with respect to the Euclidean norm):

$$\mu_{\mathcal{L}_x}^{\mathbb{A}} \geq \mu_{\mathcal{X}} \text{PWidth}(\mathcal{A})^2 \quad \text{and} \quad \mu_{\mathcal{L}_y}^{\mathbb{A}} \geq \mu_{\mathcal{Y}} \text{PWidth}(\mathcal{B})^2. \quad (72)$$

B.5 The bilinearity coefficient

In our proof, we need to relate the gradient at the point $(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})$ with the one at the point $(\mathbf{x}^*, \mathbf{y}^*)$. We can use the Lipschitz continuity of the gradient for this. We define below affine invariant quantities that can upper bound this difference.

Bilinearity coefficients. let \mathcal{L} be a strongly convex-concave function, and let $(\mathbf{x}^*, \mathbf{y}^*)$ be its unique saddle point. We define the bilinearity coefficients (M_{XY}, M_{YX}) as,

$$M_{XY} = \sup_{\substack{\mathbf{y} \in \mathcal{Y} \\ \mathbf{x}, \mathbf{s}, \mathbf{v} \in \mathcal{X} \\ \mathbf{d} = \mathbf{s} - \mathbf{v}}} \left\langle \mathbf{d}, \frac{\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{y}^*) - \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{y})}{\sqrt{\mathcal{L}^* - \mathcal{L}(\mathbf{x}^*, \mathbf{y})}} \right\rangle \quad (73)$$

and,

$$M_{YX} := \sup_{\substack{\mathbf{x} \in \mathcal{X} \\ \mathbf{y}, \mathbf{s}, \mathbf{v} \in \mathcal{Y} \\ \mathbf{d} = \mathbf{s} - \mathbf{v}}} \left\langle \mathbf{d}, \frac{\nabla_{\mathbf{y}} \mathcal{L}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{y}} \mathcal{L}(\mathbf{x}^*, \mathbf{y})}{\sqrt{\mathcal{L}(\mathbf{x}, \mathbf{y}^*) - \mathcal{L}^*}} \right\rangle. \quad (74)$$

We also define the global bilinearity coefficient as

$$M_{\mathcal{L}} := \max\{M_{XY}, M_{YX}\}. \quad (75)$$

We can upper bound these affine invariant constants with the Lipschitz constant of the gradient, the uniform strong convex-concavity constants and the diameters of the sets.

Proposition 13. *If \mathcal{X} and \mathcal{Y} are compact, $\nabla \mathcal{L}$ is Lipschitz continuous and \mathcal{L} is uniformly strongly convex-concave with constants $(\mu_{\mathcal{X}}, \mu_{\mathcal{Y}})$, then*

$$M_{XY} \leq \sqrt{\frac{2}{\mu_{\mathcal{Y}}}} L_{XY} \cdot D_{\mathcal{X}} \quad \text{and} \quad M_{YX} \leq \sqrt{\frac{2}{\mu_{\mathcal{X}}}} L_{YX} \cdot D_{\mathcal{Y}} \quad (76)$$

where L_{XY} and L_{YX} are the partial Lipschitz constants defined in Equation (40). The quantity $D_{\mathcal{X}}$ is the diameter of the compact set \mathcal{X} and $D_{\mathcal{Y}}$ is the diameter of \mathcal{Y} .

Proof.

$$\begin{aligned}
 M_{XY} &= \sup_{\substack{\mathbf{y} \in \mathcal{Y} \\ \mathbf{x}, \mathbf{s}, \mathbf{v} \in \mathcal{X} \\ \mathbf{d} = \mathbf{s} - \mathbf{v}}} \left\langle \mathbf{d}, \frac{\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{y}^*) - \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{y})}{\sqrt{\mathcal{L}^* - \mathcal{L}(\mathbf{x}^*, \mathbf{y})}} \right\rangle \\
 &\leq \sup_{\substack{\mathbf{y} \in \mathcal{Y} \\ \mathbf{x}, \mathbf{s}, \mathbf{v} \in \mathcal{X} \\ \mathbf{d} = \mathbf{s} - \mathbf{v}}} \frac{\|\mathbf{d}\|_{\mathcal{X}} \|\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{y}^*) - \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{y})\|_{\mathcal{X}^*}}{\sqrt{\mathcal{L}^* - \mathcal{L}(\mathbf{x}^*, \mathbf{y})}} \\
 &\leq \sup_{\substack{\mathbf{y} \in \mathcal{Y} \\ \mathbf{s}, \mathbf{v} \in \mathcal{X} \\ \mathbf{d} = \mathbf{s} - \mathbf{v}}} \frac{\|\mathbf{d}\|_{\mathcal{X}} L_{XY} \|\mathbf{y}^* - \mathbf{y}\|_{\mathcal{Y}}}{\sqrt{\mathcal{L}^* - \mathcal{L}(\mathbf{x}^*, \mathbf{y})}} \\
 &\leq \sup_{\mathbf{y} \in \mathcal{Y}} D_{\mathcal{X}} L_{XY} \frac{\|\mathbf{y}^* - \mathbf{y}\|_{\mathcal{Y}}}{\sqrt{\mathcal{L}^* - \mathcal{L}(\mathbf{x}^*, \mathbf{y})}}.
 \end{aligned}$$

Then using the relation between $\|\mathbf{y}^* - \mathbf{y}\|_{\mathcal{Y}}$ and $\sqrt{\mathcal{L}^* - \mathcal{L}(\mathbf{x}^*, \mathbf{y})}$ due to strong convexity (Proposition 11)

$$M_{XY} \leq \sqrt{\frac{2}{\mu_{\mathcal{Y}}}} L_{XY} \cdot D_{\mathcal{X}}. \quad (77)$$

We use a similar argument for M_{YX} which allows us to conclude. \square

B.6 Relation between the primal suboptimality

In this section, we are going to show that if the objective function \mathcal{L} is uniformly strongly convex-concave, then we have a relation between h_t and w_t . First let us introduce affine invariant constants to relate these quantities (in the context of a given saddle point $(\mathbf{x}^*, \mathbf{y}^*)$):

$$P_{\mathcal{X}} := \sup_{\mathbf{x} \in \mathcal{X}} \frac{\langle \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \hat{\mathbf{y}}(\mathbf{x})), \mathbf{x} - \mathbf{x}^* \rangle}{\sqrt{\mathcal{L}(\mathbf{x}, \mathbf{y}^*) - \mathcal{L}(\mathbf{x}^*, \mathbf{y}^*)}} \quad \text{and} \quad P_{\mathcal{Y}} := \sup_{\mathbf{y} \in \mathcal{Y}} \frac{\langle \nabla_{\mathbf{y}} \mathcal{L}(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{y}), \mathbf{y} - \mathbf{y}^* \rangle}{\sqrt{\mathcal{L}(\mathbf{x}^*, \mathbf{y}^*) - \mathcal{L}(\mathbf{x}^*, \mathbf{y})}}, \quad (78)$$

where $\hat{\mathbf{y}}(\mathbf{x}) := \arg \max_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}(\mathbf{x}, \mathbf{y})$ and $\hat{\mathbf{x}}(\mathbf{y}) := \arg \min_{\mathbf{x} \in \mathcal{X}} \mathcal{L}(\mathbf{x}, \mathbf{y})$. We also define:

$$P_{\mathcal{L}} := \max\{P_{\mathcal{X}}, P_{\mathcal{Y}}\}. \quad (79)$$

These constants can be upper bounded by easily computable constants.

Proposition 14. *For any $(\mu_{\mathcal{X}}, \mu_{\mathcal{Y}})$ -uniformly convex-concave function \mathcal{L} ,*

$$P_{\mathcal{X}} \leq \sqrt{\frac{2}{\mu_{\mathcal{X}}}} \sup_{\mathbf{z} \in \mathcal{X} \times \mathcal{Y}} \|\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{z})\|_{\mathcal{X}^*} \quad \text{and} \quad P_{\mathcal{Y}} \leq \sqrt{\frac{2}{\mu_{\mathcal{Y}}}} \sup_{\mathbf{z} \in \mathcal{X} \times \mathcal{Y}} \|\nabla_{\mathbf{y}} \mathcal{L}(\mathbf{z})\|_{\mathcal{Y}^*}. \quad (80)$$

Proof. Let us start from the definition of $P_{\mathcal{X}}$, let $\mathbf{x} \in \mathcal{X}$,

$$\begin{aligned}
 \frac{\langle \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \hat{\mathbf{y}}(\mathbf{x})), \mathbf{x} - \mathbf{x}^* \rangle}{\sqrt{\mathcal{L}(\mathbf{x}, \mathbf{y}^*) - \mathcal{L}(\mathbf{x}^*, \mathbf{y}^*)}} &\leq \frac{\|\mathbf{x} - \mathbf{x}^*\|_{\mathcal{X}} \cdot \sup(\|\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{z})\|_{\mathcal{X}^*})}{\sqrt{\mathcal{L}(\mathbf{x}, \mathbf{y}^*) - \mathcal{L}(\mathbf{x}^*, \mathbf{y}^*)}} \\
 &\leq \sqrt{\frac{2}{\mu_{\mathcal{X}}}} \sup_{\mathbf{z} \in \mathcal{X} \times \mathcal{Y}} \|\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{z})\|_{\mathcal{X}^*} \quad (\text{by strong convexity.})
 \end{aligned}$$

The same way we can get

$$P_{\mathcal{Y}} \leq \sqrt{\frac{2}{\mu_{\mathcal{Y}}}} \sup_{\mathbf{z} \in \mathcal{X} \times \mathcal{Y}} \|\nabla_{\mathbf{y}} \mathcal{L}(\mathbf{z})\|_{\mathcal{Y}^*}. \quad (81)$$

It concludes our proof. \square

One way to compute an upper bound on the supremum of the gradient is to use any reference point \bar{z} of the set:

$$\forall \bar{z} \in \mathcal{X} \times \mathcal{Y}, \quad \sup_{z \in \mathcal{X} \times \mathcal{Y}} \|\nabla_x \mathcal{L}(z)\|_{\mathcal{X}^*} \leq \nabla_x \mathcal{L}(\bar{z}) + L_{XX} D_{\mathcal{X}} + L_{XY} D_{\mathcal{Y}}. \quad (82)$$

We recall that L_{XX} is the largest (with respect to \mathbf{y}) Lipschitz constant of $\mathbf{x} \mapsto \nabla_x \mathcal{L}(\mathbf{x}, \mathbf{y})$. Note that L_{XX} is upper bounded by the global Lipschitz constant of $\nabla \mathcal{L}$. We can compute an upper bound on the supremum of the norm of $\nabla_y \mathcal{L}$ the same way.

With these above defined affine invariant constants, we can finally relate the two primal suboptimality as $h_t \leq \mathcal{O}(\sqrt{w_t})$.

Proposition 15. *For a $(\mu_{\mathcal{X}}, \mu_{\mathcal{Y}})$ -uniformly strongly convex-concave function \mathcal{L} ,*

$$h_t \leq P_{\mathcal{L}} \sqrt{2w_t} \quad \text{and} \quad P_{\mathcal{L}} \leq \sqrt{2} \sup_{z \in \mathcal{X} \times \mathcal{Y}} \left\{ \frac{\|\nabla_x \mathcal{L}(z)\|_{\mathcal{X}^*}}{\sqrt{\mu_{\mathcal{X}}}}, \frac{\|\nabla_y \mathcal{L}(z)\|_{\mathcal{Y}^*}}{\sqrt{\mu_{\mathcal{Y}}}} \right\}. \quad (83)$$

Proof. We will first work on $h_t^{(x)}$:

$$\begin{aligned} h_t^{(x)} &= \mathcal{L}(\mathbf{x}^{(t)}, \hat{\mathbf{y}}^{(t)}) - \mathcal{L}^* \\ &\leq \mathcal{L}(\mathbf{x}^{(t)}, \hat{\mathbf{y}}^{(t)}) - \mathcal{L}(\mathbf{x}^*, \hat{\mathbf{y}}^{(t)}) \\ &\leq \left\langle \mathbf{x}^{(t)} - \mathbf{x}^*, \nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \hat{\mathbf{y}}^{(t)}) \right\rangle \quad (\text{by convexity}) \\ &\leq P_{\mathcal{X}} \sqrt{w_t^{(x)}} \quad (\text{def of } P_{\mathcal{X}} \text{ (78)}). \end{aligned}$$

We can do the same thing for $h_t^{(y)}$ and $w_t^{(y)}$, thus

$$h_t \leq P_{\mathcal{L}} \left(\sqrt{w_t^{(x)}} + \sqrt{w_t^{(y)}} \right) \leq P_{\mathcal{L}} \sqrt{2w_t}, \quad (84)$$

where the last inequality uses $\sqrt{a} + \sqrt{b} \leq \sqrt{2(a+b)}$. Finally, the inequality on $P_{\mathcal{L}}$ is from Proposition 14. \square

C Relations between primal suboptimality and dual gaps

C.1 Primal suboptimality

Recall that we introduced $\hat{\mathbf{x}}^{(t)} := \arg \min_{\mathbf{x} \in \mathcal{X}} \mathcal{L}(\mathbf{x}, \mathbf{y}^{(t)})$ and similarly $\hat{\mathbf{y}}^{(t)} := \arg \max_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y})$. Then the primal suboptimality is the positive quantity

$$h_t := \mathcal{L}(\mathbf{x}^{(t)}, \hat{\mathbf{y}}^{(t)}) - \mathcal{L}(\hat{\mathbf{x}}^{(t)}, \mathbf{y}^{(t)}). \quad (85)$$

To get a convergence rate, one has to upper bound the primal suboptimality defined in (85), but it is hard to work with the moving quantities $\hat{\mathbf{x}}^{(t)}$ and $\hat{\mathbf{y}}^{(t)}$ in the analysis. This is why we use in our analysis a different merit function that uses the (fixed) saddle point $(\mathbf{x}^*, \mathbf{y}^*)$ of \mathcal{L} in its definition. We recall its definition below.

Second primal suboptimality. We define the second primal suboptimality for \mathcal{L} of the iterate $(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})$ with respect to the saddle point $(\mathbf{x}^*, \mathbf{y}^*)$ as the positive quantity:

$$w_t := \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^*) - \mathcal{L}(\mathbf{x}^*, \mathbf{y}^{(t)}). \quad (86)$$

It follows from $\mathcal{L}(\mathbf{x}^{(t)}, \hat{\mathbf{y}}^{(t)}) \geq \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^*)$ and $\mathcal{L}(\mathbf{x}^*, \mathbf{y}^{(t)}) \geq \mathcal{L}(\hat{\mathbf{x}}^{(t)}, \mathbf{y}^{(t)})$ that $w_t \leq h_t$. Furthermore, under the assumption of uniform strong convex-concavity, we proved in Proposition 15 that the square root of w_t upper bounds h_t up to a constant.

C.2 Gap inequalities

In this section, we will prove the crucial inequalities relating suboptimality and the gap function. Let's recall the definition of $\mathbf{s}^{(t)}$ and $\mathbf{v}^{(t)}$:

$$\mathbf{s}^{(t)} := \arg \min_{\mathbf{s} \in \mathcal{X} \times \mathcal{Y}} \langle \mathbf{s}, \mathbf{r}^{(t)} \rangle \quad \text{and} \quad \mathbf{v}^{(t)} := \arg \max_{\mathbf{v} \in \mathcal{S}_x^t \times \mathcal{S}_y^t} \langle \mathbf{v}, \mathbf{r}^{(t)} \rangle \quad (87)$$

where $(\mathbf{r}^{(t)})^\top := ((\mathbf{r}_x^{(t)})^\top, (\mathbf{r}_y^{(t)})^\top) := (\nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}), -\nabla_y \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}))$. Also, the following various gaps are defined as

$$g_t^{\text{FW}} := \langle \mathbf{d}_{\text{FW}}^{(t)}, -\mathbf{r}^{(t)} \rangle, \quad g_t^{\text{PFW}} := \langle \mathbf{d}_{\text{PFW}}^{(t)}, -\mathbf{r}^{(t)} \rangle \quad \text{and} \quad g_t := \langle \mathbf{d}^{(t)}, -\mathbf{r}^{(t)} \rangle \quad (88)$$

where $\mathbf{d}_{\text{FW}}^{(t)} = \mathbf{s}^{(t)} - \mathbf{z}^{(t)}$ and $\mathbf{d}_{\text{PFW}}^{(t)} = \mathbf{s}^{(t)} - \mathbf{v}^{(t)}$. The direction $\mathbf{d}^{(t)}$ is the direction chosen by the algorithm at step t : it is always $\mathbf{d}_{\text{FW}}^{(t)}$ for SP-FW, and can be either $\mathbf{d}_{\text{FW}}^{(t)}$ or $\mathbf{d}_{\text{A}}^{(t)} := \mathbf{z}^{(t)} - \mathbf{v}^{(t)}$ for SP-AFW. Even if the definitions of these gaps are different, the formalism for the analysis of the convergence of both algorithms is going to be fairly similar. It is straightforward to notice that $g_t^{\text{PFW}} \geq g_t$ and one can show that the current gap g_t is lower bounded by half of g_t^{PFW} :

Lemma 16. *For the SP-AFW algorithm, the current gap g_t can be bounded as follows:*

$$\frac{1}{2} g_t^{\text{PFW}} \leq g_t \leq g_t^{\text{PFW}} \quad (89)$$

Proof. First let's show the RHS of the inequality,

$$g_t^{\text{PFW}} := \langle \mathbf{d}_{\text{PFW}}^{(t)}, -\mathbf{r}^{(t)} \rangle = \langle \mathbf{d}_{\text{A}}^{(t)}, -\mathbf{r}^{(t)} \rangle + \langle \mathbf{d}_{\text{FW}}^{(t)}, -\mathbf{r}^{(t)} \rangle \geq \langle \mathbf{d}^{(t)}, -\mathbf{r}^{(t)} \rangle \quad (90)$$

because both $\langle \mathbf{d}_{\text{A}}^{(t)}, -\mathbf{r}^{(t)} \rangle \geq 0$ and $\langle \mathbf{d}_{\text{FW}}^{(t)}, -\mathbf{r}^{(t)} \rangle \geq 0$ from their definition. For the LHS inequality, we use the fact that $g_t = \max \left\{ \langle \mathbf{d}_{\text{A}}^{(t)}, -\mathbf{r}^{(t)} \rangle, \langle \mathbf{d}_{\text{FW}}^{(t)}, -\mathbf{r}^{(t)} \rangle \right\}$ for SP-AFW and thus:

$$g_t^{\text{PFW}} = \langle \mathbf{d}_{\text{A}}^{(t)}, -\mathbf{r}^{(t)} \rangle + \langle \mathbf{d}_{\text{FW}}^{(t)}, -\mathbf{r}^{(t)} \rangle \leq 2g_t. \quad (91)$$

□

In the following, we will assume that we are in one of the two following cases:

- (I) The saddle point of \mathcal{L} belongs to the relative interior of $\mathcal{X} \times \mathcal{Y}$.
- (P) \mathcal{X} and \mathcal{Y} are polytopes, i.e. $\exists \mathcal{A}, \mathcal{B}$ finite s.t. $\mathcal{X} = \text{conv}(\mathcal{A})$, $\mathcal{Y} = \text{conv}(\mathcal{B})$.

Then either $\mu_{\mathcal{L}}^{\text{int}} > 0$ (case I) or $\mu_{\mathcal{L}}^{\text{A}} > 0$ (case P). Let's write the gap function as the sum of two smaller gap functions:

$$g_t = \underbrace{\langle \mathbf{d}_{(x)}^{(t)}, -\mathbf{r}_x^{(t)} \rangle}_{=: g_t^{(x)}} + \underbrace{\langle \mathbf{d}_{(y)}^{(t)}, -\mathbf{r}_y^{(t)} \rangle}_{=: g_t^{(y)}} \quad (92)$$

Because of the convex-concavity of \mathcal{L} , this scalar product bounds the differences between the value of \mathcal{L} at the point $(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})$ and the value of \mathcal{L} at another point. Hence this gap function upper-bounds h_t and w_t defined in (85) and (86). More concretely, we have the following lemma.

Lemma 17. For all t in \mathbb{N} , $\mathbf{x} \in \mathcal{X}$ and $\mathbf{y} \in \mathcal{Y}$

$$g_t^{\text{PFW}} \geq g_t^{\text{FW}} \geq \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}) - \mathcal{L}(\mathbf{x}, \mathbf{y}^{(t)}), \quad (93)$$

and, furthermore,

$$g_t \geq h_t \geq w_t. \quad (94)$$

Proof. First let's show the LHS of (93),

$$g_t^{\text{PFW}} = \left\langle \mathbf{d}_{\text{PFW}}^{(t)}, -\mathbf{r}^{(t)} \right\rangle = \left\langle \mathbf{d}_A^{(t)}, -\mathbf{r}^{(t)} \right\rangle + \left\langle \mathbf{d}_{\text{FW}}^{(t)}, -\mathbf{r}^{(t)} \right\rangle \geq \left\langle \mathbf{d}_{\text{FW}}^{(t)}, -\mathbf{r}^{(t)} \right\rangle = g_t^{\text{FW}} \quad (95)$$

because one can easily derive that $\left\langle \mathbf{d}_A^{(t)}, -\mathbf{r}^{(t)} \right\rangle \geq 0$ from the definition of the away direction $\mathbf{d}_A^{(t)}$. It follows from convexity of $\mathbf{x} \mapsto \mathcal{L}(\mathbf{x}, \mathbf{y}^{(t)})$ that for all \mathbf{x} in \mathcal{X} ,

$$(g_t^{\text{FW}})_x := \left\langle (\mathbf{d}_{\text{FW}}^{(t)})_x, -\nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \right\rangle \geq \left\langle \mathbf{x} - \mathbf{x}^{(t)}, -\nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \right\rangle \quad (96)$$

$$\geq \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}, \mathbf{y}^{(t)}). \quad (97)$$

A similar inequality emerges through the convexity of $\mathbf{y} \mapsto -\mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y})$,

$$(g_t^{\text{FW}})_y := \left\langle (\mathbf{d}_{\text{FW}}^{(t)})_y, \nabla_y \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \right\rangle \geq \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}), \quad (98)$$

which gives us

$$g_t^{\text{FW}} \geq \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}, \mathbf{y}^{(t)}) + \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}), \quad (99)$$

which shows (93). By using $\mathbf{x} = \hat{\mathbf{x}}^{(t)}$ and $\mathbf{y} = \hat{\mathbf{y}}^{(t)}$ in (93), we get $g_t^{\text{FW}} \geq h_t$. We also know that $g_t = \max(g_t^A, g_t^{\text{FW}}) \geq g_t^{\text{FW}}$ for SP-AFW. So combining with $h_t \geq w_t$ that we already knew, we get (94). \square

Next, we recall two lemmas, one from Lacoste-Julien and Jaggi (2013) and the other one from (Lacoste-Julien and Jaggi, 2015). These lemmas upper bound the primal suboptimality with the square of the gap times a constant depending on the geometric (or the interior) strong convexity constant.

Lemma 18 (Lacoste-Julien and Jaggi (2015), Lacoste-Julien and Jaggi (2013)). *If f is strongly convex, then for any $\mathbf{x}^{(t)} \in \mathcal{X}$,*

$$f(\mathbf{x}^{(t)}) - f(\mathbf{x}_c) \leq \frac{(g_t^{\text{FW}})^2}{2\mu_f^{\mathbf{x}_c}} \quad \text{if } \mathbf{x}_c \in \text{interior of } \mathcal{X} \quad (\text{Lacoste-Julien and Jaggi, 2013}) \quad (100)$$

and

$$f(\mathbf{x}^{(t)}) - f^* \leq \frac{(g_t^{\text{PFW}})^2}{2\mu_f^A} \quad \text{if } \mathcal{X} = \text{conv}(A) \quad (\text{Lacoste-Julien and Jaggi, 2015}) \quad (101)$$

where $g_t^{\text{FW}} = \langle \mathbf{x}^{(t)} - \mathbf{s}^{(t)}, \nabla f(\mathbf{x}^{(t)}) \rangle$, $g_t^{\text{PFW}} = \langle \mathbf{v}^{(t)} - \mathbf{s}^{(t)}, \nabla f(\mathbf{x}^{(t)}) \rangle$ and $f^* = \min_{\mathbf{x} \in \mathcal{X}} f$.

Notice once again that in this lemma we do not need \mathbf{x}_c to be optimal.

Proof. Let $\mathbf{x}^{(t)} \neq \mathbf{x}_c$. Using the definition of interior strong convexity (46) and choosing γ such that $\mathbf{x}_c = \mathbf{x}^{(t)} + \gamma (\bar{\mathbf{s}}(\mathbf{x}^{(t)}, \mathbf{x}_c) - \mathbf{x}^{(t)})$, we get

$$\begin{aligned} f(\mathbf{x}_c) - f(\mathbf{x}^{(t)}) &\geq \gamma \left\langle \bar{\mathbf{s}}(\mathbf{x}_c, \mathbf{x}^{(t)}) - \mathbf{x}^{(t)}, \nabla f(\mathbf{x}^{(t)}) \right\rangle + \gamma^2 \frac{\mu_f^{\mathbf{x}_c}}{2} \\ &\geq -\gamma g_t^{\text{FW}} + \gamma^2 \frac{\mu_f^{\mathbf{x}_c}}{2} \\ &\geq \frac{(g_t^{\text{FW}})^2}{2\mu_f^{\mathbf{x}_c}}. \end{aligned}$$

The last line of this derivation is obtained through the inequality: $-a^2 + 2ab - b^2 \leq 0$. If $\mathbf{x}^{(t)} = \mathbf{x}_c$ the inequality is just the positivity of the gap.

For the second statement, we will use the definition of the geometric strong convexity constant (Equation (53)) at the point $\mathbf{x} = \mathbf{x}^{(t)}$ and $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$. Recall that $\gamma^A(\mathbf{x}, \mathbf{x}^*) = \frac{\langle -\nabla f(\mathbf{x}), \mathbf{x}^* - \mathbf{x} \rangle}{\langle -\nabla f(\mathbf{x}), \mathbf{s}_f(\mathbf{x}) - \mathbf{v}_f(\mathbf{x}) \rangle}$.

$$\begin{aligned}
 f(\mathbf{x}^*) - f(\mathbf{x}^{(t)}) &\geq \left\langle \mathbf{x}^* - \mathbf{x}^{(t)}, \nabla f(\mathbf{x}^{(t)}) \right\rangle + \frac{\mu_f^A}{2} \gamma^A(\mathbf{x}^{(t)}, \mathbf{x}^*)^2 \\
 &= -\gamma^A(\mathbf{x}^{(t)}, \mathbf{x}^*) \left\langle \mathbf{s}_f(\mathbf{x}^{(t)}) - \mathbf{v}_f(\mathbf{x}^{(t)}), -\nabla f(\mathbf{x}^{(t)}) \right\rangle + \frac{\mu_f^A}{2} \gamma^A(\mathbf{x}^{(t)}, \mathbf{x}^*)^2 \\
 &\geq -\gamma^A(\mathbf{x}^{(t)}, \mathbf{x}^*) g_t^{\text{PFW}} + \frac{\mu_f^A}{2} \gamma^A(\mathbf{x}^{(t)}, \mathbf{x}^*)^2 \quad (\text{Equation (54)}) \\
 &\geq \frac{-(g_t^{\text{PFW}})^2}{2\mu_f^A}.
 \end{aligned}$$

□

Lemma 18 is useful to understand the following lemma and its proof which is just an extension to the convex-concave case.

Lemma 19 (Quadratic gap upper bound on second suboptimality for (I) or (P)). *If \mathcal{L} is a strongly convex-concave function, then for any $(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \in \mathcal{X} \times \mathcal{Y}$,*

$$w_t \leq \frac{(g_t^{\text{FW}})^2}{2\mu_{\mathcal{L}}^{\text{int}}} \quad \text{for (I)} \quad \text{and} \quad w_t \leq h_t \leq \frac{(g_t^{\text{PFW}})^2}{2\mu_{\mathcal{L}}^A} \quad \text{for (P)} \quad (102)$$

where the gaps are defined in (88), $\mu_{\mathcal{L}}^{\text{int}} := \min\{\mu_{\mathcal{L}}^{\mathbf{x}^*}, \mu_{\mathcal{L}}^{\mathbf{y}^*}\}$ (i.e. using the reference points $(\mathbf{x}_c, \mathbf{y}_c) := (\mathbf{x}^*, \mathbf{y}^*)$ in the definition (66)) and $\mu_{\mathcal{L}}^A$ is the geometric strong convex-concavity of \mathcal{L} over $\mathcal{A} \times \mathcal{B}$, as defined in (71).

Proof. For (I):

Let the function f on \mathcal{X} be defined by $f(\mathbf{x}) = \mathcal{L}(\mathbf{x}', \mathbf{y}^{(t)})$, and the function g on \mathcal{Y} be $g(\mathbf{y}') = -\mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}')$. Then using the Lemma 18 on the function f with the reference point \mathbf{x}^* , and on g with reference point \mathbf{y}^* , we get

$$\begin{aligned}
 \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^*, \mathbf{y}^{(t)}) &\leq \frac{\left\langle \mathbf{s}_x^{(t)} - \mathbf{x}^{(t)}, -\nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \right\rangle^2}{2\mu_f^{\mathbf{x}^*}} \\
 \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^*) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) &\leq \frac{\left\langle \mathbf{s}_y^{(t)} - \mathbf{y}^{(t)}, \nabla_y \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \right\rangle^2}{2\mu_g^{\mathbf{y}^*}}.
 \end{aligned}$$

As $\mu_{\mathcal{L}}^{\text{int}}$ is smaller than both $\mu_f^{\mathbf{x}^*}$ and $\mu_g^{\mathbf{y}^*}$ by the definition (66), we can use it in the denominator of the above two inequalities. As we saw from Section C.2 in (92), the gap can be split as sum of the gap of the block \mathcal{X} and the gap of the block \mathcal{Y} , i.e. $g_t^{\text{FW}} = \left\langle \mathbf{s}_x^{(t)} - \mathbf{x}^{(t)}, -\nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \right\rangle + \left\langle \mathbf{s}_y^{(t)} - \mathbf{y}^{(t)}, \nabla_y \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \right\rangle$. Then, using the inequality: $a^2 + b^2 \leq (a + b)^2$ for $(a, b \geq 0)$, we obtain

$$w_t \leq \frac{(g_t^{\text{FW}})^2}{2\mu_{\mathcal{L}}^{\text{int}}}. \quad (103)$$

For (P):

Using the Lemma 18 for case (P) on the same functions f and g defined above, we get

$$\begin{aligned}\mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}(\widehat{\mathbf{x}}^{(t)}, \mathbf{y}^{(t)}) &\leq \frac{\left\langle \mathbf{s}_x^{(t)} - \mathbf{v}_x^{(t)}, \nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \right\rangle^2}{2\mu_f^A} \\ \mathcal{L}(\mathbf{x}^{(t)}, \widehat{\mathbf{y}}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) &\leq \frac{\left\langle \mathbf{s}_y^{(t)} - \mathbf{v}_y^{(t)}, -\nabla_y \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \right\rangle^2}{2\mu_g^A}.\end{aligned}$$

Using a similar argument as the one to get (103), using that $\mu_{\mathcal{L}}^A$ is smaller than both μ_f^A and μ_g^A , and referring to the separation of the gap (92), we get

$$h_t \leq \frac{(g_t^{\text{PFW}})^2}{2\mu_{\mathcal{L}}^A}. \quad (104)$$

□

D Convergence analysis

In this section, we are going to show two important lemmas. The first one shows that under some assumptions we can get a Frank-Wolfe-style induction scheme relating the second suboptimality of the potential update w_γ , the current value of the second suboptimality w_t , the gap g_t and any step size $\gamma \in [0, \gamma_{\max}]$. The second lemma will relate the gap and the square root of w_t ; this relation enables us to get a rate on the gap after getting a rate on w_t .

D.1 First lemmas

The first lemma in this section is inspired from the standard FW progress lemma, such as Lemma C.2 in (Lacoste-Julien et al., 2013), though it requires a non-trivial change due to the compensation phenomenon for \mathcal{L} mentioned in the main text in (10). In the following, we define the possible updated iterate \mathbf{z}_γ for $\gamma \in [0, \gamma_{\max}]$:

$$\mathbf{z}_\gamma := (\mathbf{x}_\gamma, \mathbf{y}_\gamma) := \mathbf{z}^{(t)} + \gamma \mathbf{d}^{(t)}, \quad \text{where } \mathbf{d}^{(t)} \text{ is the direction of the step.} \quad (105)$$

For a FW step $\mathbf{d}^{(t)} = \mathbf{d}_{\text{FW}}^{(t)} := \mathbf{s}^{(t)} - \mathbf{z}^{(t)}$ and for an away step $\mathbf{d}^{(t)} = \mathbf{d}_{\text{A}}^{(t)} := \mathbf{z}^{(t)} - \mathbf{v}^{(t)}$. We also define the corresponding new suboptimality for \mathbf{z}_γ :

$$w_\gamma := \mathcal{L}(\mathbf{x}_\gamma, \mathbf{y}^*) - \mathcal{L}(\mathbf{x}^*, \mathbf{y}_\gamma). \quad (106)$$

Lemma 20 (Suboptimality progress for SP-FW and SP-AFW). *Let \mathcal{L} be strongly convex-concave,*

If we are in case (I) and $\mathbf{d}^{(t)} = \mathbf{d}_{\text{FW}}^{(t)}$ is a FW direction, we have for any $\gamma \in [0, 1]$:

$$w_\gamma \leq w_t - \nu^{\text{FW}} \gamma g_t^{\text{FW}} + \gamma^2 C_{\mathcal{L}}, \quad \text{where } \nu^{\text{FW}} := 1 - \frac{M_{\mathcal{L}}}{\sqrt{\mu_{\mathcal{L}}^{\text{int}}}}. \quad (107)$$

If we are in case (P) and $\mathbf{d}^{(t)}$ is defined from a step of SP-AFW (Algorithm 3), we have for any $\gamma \in [0, \gamma_{\max}]$:

$$w_\gamma \leq w_t - \nu^{\text{PFW}} \gamma g_t^{\text{PFW}} + \gamma^2 C_{\mathcal{L}}, \quad \text{where } \nu^{\text{PFW}} := \frac{1}{2} - \frac{M_{\mathcal{L}}}{\sqrt{\mu_{\mathcal{L}}^A}}. \quad (108)$$

Proof. The beginning of the argument works with any direction $\mathbf{d}^{(t)}$. Recall that $w_t^{(x)} := \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^*) - \mathcal{L}^*$ and $w_t^{(y)} := -\mathcal{L}(\mathbf{x}^*, \mathbf{y}^{(t)}) + \mathcal{L}^*$. Now writing $\mathbf{x}_\gamma = \mathbf{x}^{(t)} + \gamma \mathbf{d}_x^{(t)}$ and using the definition of the curvature C_f (41) for the function $\mathbf{x} \mapsto f(\mathbf{x}) := \mathcal{L}(\mathbf{x}, \mathbf{y}^*)$, we get

$$\begin{aligned} w_\gamma^{(x)} &:= \mathcal{L}(\mathbf{x}_\gamma, \mathbf{y}^*) - \mathcal{L}^* \\ &\leq \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^*) - \mathcal{L}^* + \gamma \left\langle \mathbf{d}_x^{(t)}, \nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^*) \right\rangle + \gamma^2 \frac{C_{\mathcal{L}}}{2}, \end{aligned} \quad (109)$$

since $C_{\mathcal{L}} \leq C_f$ by definition (58). Recall that the gap function g_t can be decomposed by (92) into two smaller gap functions $g_t^{(x)} := \left\langle \mathbf{d}_x^{(t)}, -\mathbf{r}_x^{(t)} \right\rangle$ and $g_t^{(y)} := \left\langle \mathbf{d}_y^{(t)}, -\mathbf{r}_y^{(t)} \right\rangle$. We define $\epsilon_t := \left\langle \mathbf{d}_x^{(t)}, \nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^*) - \nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \right\rangle$ to be the sequence representing the error between the gradient used for the minimization and the gradient at the point $(\mathbf{x}^{(t)}, \mathbf{y}^*)$. Then,

$$w_\gamma^{(x)} \leq w_t^{(x)} - \gamma g_t^{(x)} + \gamma \epsilon_t + \gamma^2 \frac{C_{\mathcal{L}}}{2}. \quad (110)$$

Now, as M_{XY} is finite (under Lipschitz gradient assumption), we can use the definition of the bilinearity constant (73) to get

$$|\epsilon_t| = \left| \left\langle \mathbf{d}_x^{(t)}, \nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^*) - \nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \right\rangle \right| \leq \sqrt{w_t^{(y)}} M_{XY}. \quad (111)$$

Combining equations (110) and (111) we finally obtain

$$w_\gamma^{(x)} \leq w_t^{(x)} - \gamma g_t^{(x)} + \gamma M_{XY} \sqrt{w_t^{(y)}} + \gamma^2 \frac{C_{\mathcal{L}}}{2}. \quad (112)$$

We can get an analogous inequality for $w_\gamma^{(y)}$,

$$w_\gamma^{(y)} \leq w_t^{(y)} - \gamma g_t^{(y)} + \gamma M_{YX} \sqrt{w_t^{(x)}} + \gamma^2 \frac{C_{\mathcal{L}}}{2}. \quad (113)$$

Then adding $w_\gamma^{(x)}$ and $w_\gamma^{(y)}$ and using $\sqrt{a} + \sqrt{b} \leq \sqrt{2(a+b)}$ (coming from the concavity of $\sqrt{\cdot}$), we get

$$w_\gamma \leq w_t - \gamma g_t + \gamma M_{\mathcal{L}} \sqrt{2w_t} + 2\gamma^2 \frac{C_{\mathcal{L}}}{2}. \quad (114)$$

We stress that the above inequality (114) is valid for any direction $\mathbf{d}^{(t)}$, using $g_t := \langle \mathbf{d}^{(t)}, -\mathbf{r}^{(t)} \rangle$, and for any feasible step size γ such that $\mathbf{z}_\gamma \in \mathcal{X} \times \mathcal{Y}$ (the last condition was used in the definition of C_f ; see also footnote 2 for more information).

To finish the argument, we now use the specific property of the direction $\mathbf{d}^{(t)}$ and use the crucial Lemma 19 that relates w_t with the square of the appropriate gap.

For the case (I) of interior saddle point, we consider $\mathbf{d}^{(t)} = \mathbf{d}_{\text{FW}}^{(t)}$ and thus $g_t = g_t^{\text{FW}}$. Then combining Lemma 19 (using the interior strong convexity constant) with (114), we get

$$w_\gamma \leq w_t - \gamma \left(1 - \frac{M_{\mathcal{L}}}{\sqrt{\mu_{\mathcal{L}}^{\text{int}}}} \right) g_t^{\text{FW}} + \gamma^2 C_{\mathcal{L}}. \quad (115)$$

For the case (P) of polytope domains, we consider $\mathbf{d}^{(t)}$ as defined by the SP-AFW algorithm. We thus have $g_t \geq \frac{1}{2} g_t^{\text{PFW}}$ by Lemma 16. Then combining Lemma 19 (using the geometric strong convexity constant) with (114), we get

$$w_\gamma \leq w_t - \gamma \left(\frac{1}{2} - \frac{M_{\mathcal{L}}}{\sqrt{\mu_{\mathcal{L}}^{\text{A}}}} \right) g_t^{\text{PFW}} + \gamma^2 C_{\mathcal{L}}, \quad (116)$$

which finishes the proof. Still in the case (P), we also present an inequality in terms of the direction gap g_t (which yields a better constant that will be important for the sublinear convergence proof in Theorem 25) by using instead the inequality $g_t^{\text{PFW}} \geq g_t$ (Lemma 16) with Lemma 19 and (114):

$$w_\gamma \leq w_t - \gamma \left(1 - \frac{M_{\mathcal{L}}}{\sqrt{\mu_{\mathcal{L}}^A}} \right) g_t + \gamma^2 C_{\mathcal{L}}. \quad (117)$$

□

The above lemma uses a specific update direction $\mathbf{d}^{(t)}$ to get a potential new suboptimality w_γ . By using the property that $w_\gamma \geq 0$ always, we can actually derive an upper bound on the gap in terms of w_t irrespective of any algorithm (i.e. this relationship holds for any possible feasible point $(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})$). More precisely, for SP-FW algorithm (case (I)) the only thing we need to set is a feasible point $\mathbf{z}^{(t)}$ but for the SP-AFW algorithm (case (P)) we also need an active set expansion for $\mathbf{z}^{(t)}$ for which the maximum away step size is larger than $\frac{\nu g_t}{2C_{\mathcal{L}}}$ (which can potentially not be the active set calculated by an algorithm). This is stated in the following theorem, which is a saddle point generalization of the gap upper bound given in Theorem 2 of (Lacoste-Julien and Jaggi, 2015).

Theorem 21 (Bounding the gap with the second suboptimality). *If \mathcal{L} is strongly convex-concave and has a finite curvature constant then*

- case (I): For any $\mathbf{z}^{(t)} \in \mathcal{X} \times \mathcal{Y}$,

$$g_t^{\text{FW}} \leq \frac{2}{\nu^{\text{FW}}} \max \left\{ \sqrt{C_{\mathcal{L}} w_t}, w_t \right\}. \quad (118)$$

Since $\mathbf{z}^{(t)}$ is fixed, this statement is algorithm free.

- case (P): For any $\mathbf{z}^{(t)} \in \mathcal{X} \times \mathcal{Y}$, if there exists an active set expansion for $\mathbf{z}^{(t)}$ for which $\gamma_{\max} = 1$ or $\gamma_{\max} \geq \frac{\nu^{\text{PFW}} g_t^{\text{PFW}}}{2C_{\mathcal{L}}}$ (see (122) for the definition of γ_{\max}) then,

$$g_t^{\text{PFW}} \leq \frac{2}{\nu^{\text{PFW}}} \max \left\{ \sqrt{C_{\mathcal{L}} w_t}, w_t \right\}. \quad (119)$$

Both statement are algorithm free but g_t^{PFW} depends on a chosen expansion of $\mathbf{z}^{(t)}$:

$$\mathbf{z}^{(t)} = \sum_{(\mathbf{v}_x, \mathbf{v}_y) \in \mathcal{S}^{(t)}} \alpha_{\mathbf{v}_x}^{(t)} \alpha_{\mathbf{v}_y}^{(t)} (\mathbf{v}_x, \mathbf{v}_y) \quad \text{where} \quad \mathcal{S}^{(t)} := \left\{ (\mathbf{v}_x, \mathbf{v}_y) \in \mathcal{A} \times \mathcal{B} ; \alpha_{\mathbf{v}_x}^{(t)} \alpha_{\mathbf{v}_y}^{(t)} > 0 \right\}, \quad (120)$$

because,

$$\begin{aligned} g_t^{\text{PFW}} &:= \langle -\mathbf{r}^{(t)}, \mathbf{d}_{\text{FW}}^{(t)} + \mathbf{d}_{\text{A}}^{(t)} \rangle \\ \text{where} \quad \mathbf{d}_{\text{A}}^{(t)} &:= \mathbf{z}^{(t)} - \arg \max_{\mathbf{v} \in \mathcal{S}_x \times \mathcal{S}_y} \langle \mathbf{r}^{(t)}, \mathbf{v} \rangle, \\ \text{and} \quad \mathbf{d}_{\text{FW}}^{(t)} &:= \arg \min_{\mathbf{v} \in \mathcal{A} \times \mathcal{B}} \langle \mathbf{r}^{(t)}, \mathbf{v} \rangle - \mathbf{z}^{(t)}. \end{aligned} \quad (121)$$

The maximum step size associated with the active set expansion described in Equation (120) is

$$\gamma_{\max} := \begin{cases} 1 & \text{if } \langle -\mathbf{r}^{(t)}, \mathbf{d}_{\text{A}}^{(t)} \rangle \leq \langle -\mathbf{r}^{(t)}, \mathbf{d}_{\text{FW}}^{(t)} \rangle, \\ \min \left\{ \frac{\alpha_{\mathbf{v}_x}^{(t)}}{1 - \alpha_{\mathbf{v}_x}^{(t)}}, \frac{\alpha_{\mathbf{v}_y}^{(t)}}{1 - \alpha_{\mathbf{v}_y}^{(t)}} \right\} & \text{otherwise.} \end{cases} \quad (122)$$

Proof. In this proof, we let (g_t, ν) to stand respectively for $(g_t^{\text{FW}}, \nu^{\text{FW}})$ for case (I) or $(g_t^{\text{PFW}}, \nu^{\text{PFW}})$ for case (P). We will start from the inequalities (107) and (108) in Lemma 20. Equation (107) is valid considering a FW direction $\mathbf{d}_{\text{FW}}^{(t)}$, Equation (108) is valid if we consider the direction that would have been set by the SP-AFW algorithm if it was run at point $\mathbf{z}^{(t)}$ with the active set expansion described in the theorem statement. Since $w_\gamma \geq 0$, for both cases become :

$$0 \leq w_t - \gamma \nu g_t + \gamma^2 C_{\mathcal{L}}, \quad (123)$$

then we can put the gap on the LHS,

$$\gamma \nu g_t - \gamma^2 C_{\mathcal{L}} \leq w_t. \quad (124)$$

This inequality is valid for any $\gamma \in [0, \gamma_{\max}]$. In order to get the tightest bound between the gap and the suboptimality, we will maximize the LHS. It can be maximized with $\bar{\gamma} := \frac{\nu g_t}{2C_{\mathcal{L}}} = \gamma_t$. Now we have two cases:

If $\bar{\gamma} \in [0, \gamma_{\max}]$, then we get: $\nu g_t \leq 2\sqrt{C_{\mathcal{L}} w_t}$.

And if $\gamma_{\max} = 1$ and $\bar{\gamma} = \frac{\nu g_t}{2C_{\mathcal{L}}} > 1$, then setting $\gamma = 1$ we get: $\nu g_t \leq 2w_t$. By taking the maximum between the two options, we get the theorem statement. \square

The previous theorem guarantees that the gap gets small when w_t gets small (only for the non-drop steps if in situation (P)). As we use the gap as a stopping criterion in the algorithm, this is a useful theorem to provide an upper bound on the number of iterations needed to get a certificate of suboptimality.

The following corollary provides a better bound on h_t than the inequality $h_t \leq \text{cst}\sqrt{w_t}$ previously shown (14) when the situation is (P). It will be useful later to get a better rate of convergence for h_t under hypothesis (P).

Corollary 22 (Tighter bound on h_t for non-drop steps in situation (P)). *Suppose that \mathcal{L} is strongly convex-concave and has a finite curvature constant, and that the domain is a product of polytopes (i.e. we are in situation (P)). Let $\mathbf{z}^{(t)} \in \mathcal{X} \times \mathcal{Y}$ be given. If there exists an active set expansion for $\mathbf{z}^{(t)}$ for which the maximum step size is larger than $\frac{\nu g_t}{2C_{\mathcal{L}}}$ (see Theorem (21) for more details),*

$$h_t \leq \frac{2 \max\{C_{\mathcal{L}}, w_t\}}{\mu_{\mathcal{L}}^{\text{A}} (\nu^{\text{PFW}})^2} w_t, \quad (125)$$

where ν^{PFW} is defined in (108).

Proof. By Lemma 19 in situation (P), we have:

$$h_t \leq \frac{(g_t^{\text{PFW}})^2}{2\mu_{\mathcal{L}}^{\text{A}}} \leq \frac{2 \max\{C_{\mathcal{L}}, w_t\}}{\mu_{\mathcal{L}}^{\text{A}} (\nu^{\text{PFW}})^2} w_t.$$

The last inequality is obtained by applying the upper bound on the gap given in the previous Theorem 21. \square

D.2 Proof of Theorem 1

In this section, we will prove that under some conditions on the constant defined in subsections B.2 and B.3, the suboptimalities w_t vanish linearly with the adaptive step size $\gamma_t = \min \left\{ \gamma_{\max}, \frac{\nu}{2C_{\mathcal{L}}} g_t \right\}$ or sublinearly with the universal step size $\gamma_t = \min \left\{ \gamma_{\max}, \frac{2}{2+k(t)} \right\}$.

Lemma 23 (Geometric decrease of second suboptimality). *Let \mathcal{L} be a strongly convex-concave function with a smoothness constant $C_{\mathcal{L}}$, a positive interior strong convex-concavity constant $\mu_{\mathcal{L}}^{\text{int}}$ (66) or a positive geometric strong convex-concavity $\mu_{\mathcal{L}}^{\text{A}}$ (71). Let us also define the rate multipliers ν as*

$$\nu^{\text{FW}} := 1 - \frac{M_{\mathcal{L}}}{\sqrt{\mu_{\mathcal{L}}^{\text{int}}}} \quad \text{and} \quad \nu^{\text{PFW}} := \frac{1}{2} - \frac{M_{\mathcal{L}}}{\sqrt{\mu_{\mathcal{L}}^{\text{A}}}} \quad (\text{see Equation (75) for the definition of } M_{\mathcal{L}}). \quad (126)$$

Let the tuple $(g_t, \nu, \mu_{\mathcal{L}})$ refers to either $(g_t^{\text{FW}}, \nu^{\text{FW}}, \mu_{\mathcal{L}}^{\text{int}})$ for case (I) where the algorithm is SP-FW, or $(g_t^{\text{PFW}}, \nu^{\text{PFW}}, \mu_{\mathcal{L}}^{\text{A}})$ for case (P) where the algorithm is SP-AFW,

If $\nu > 0$, then at each non-drop step (when $\gamma_t < \gamma_{\max}$ or $\gamma_{\max} \geq 1$), the suboptimality w_t of the algorithm with step size $\gamma_t = \min(\gamma_{\max}, \frac{\nu}{2C_{\mathcal{L}}}g_t)$ decreases geometrically as

$$w_{t+1} \leq (1 - \rho_{\mathcal{L}})w_t \quad (127)$$

where $\rho_{\mathcal{L}} := \frac{\nu^2 \mu_{\mathcal{L}}}{2C_{\mathcal{L}}}$. Moreover, for case (I) there is no drop step and for case (P) the number of drop step (when $\gamma_t = \gamma_{\max}$) is upper bounded by two third of the number of iteration (see Section A, Equation (38)), while when we have a drop step, we still have:

$$w_{t+1} \leq w_t. \quad (128)$$

Proof. The bulk of the proof is of a similar form for both SP-FW and SP-AFW, and so in the following, we let $(g_t, \mu_{\mathcal{L}}, \nu)$ to stand respectively for $(g_t^{\text{FW}}, \mu_{\mathcal{L}}^{\text{int}}, \nu^{\text{FW}})$ for SP-FW (case (I)) or $(g_t^{\text{PFW}}, \mu_{\mathcal{L}}^{\text{A}}, \nu^{\text{PFW}})$ for SP-AFW (case (P)). As $\gamma_t \leq \gamma_{\max}$, we can apply the important Lemma 20 with $\gamma = \gamma_t$ (the actual step size that was taken in the algorithm) to get:

$$w_{t+1} = w_{\gamma_t} \leq w_t - \nu\gamma_t g_t + \gamma_t^2 C_{\mathcal{L}}. \quad (129)$$

We note in passing that the adaptive step size rule $\gamma_t = \min(\gamma_{\max}, \frac{\nu}{2C_{\mathcal{L}}}g_t)$ was specifically chosen to minimize the RHS of (129) among the feasible step sizes.

If $\frac{\nu}{2C_{\mathcal{L}}}g_t \leq \gamma_{\max}$, then we have $\gamma_t = \frac{\nu}{2C_{\mathcal{L}}}g_t$ and so (129) becomes:

$$w_{t+1} \leq w_t - \frac{\nu^2}{2C_{\mathcal{L}}}(g_t)^2 + \frac{\nu^2}{4C_{\mathcal{L}}}(g_t)^2 = w_t - \frac{\nu^2}{4C_{\mathcal{L}}}(g_t)^2. \quad (130)$$

Applying the fact that the square of the appropriate gap upper bounds w_t (Lemma 19 with a similar form for both cases (I) and (P)), we directly obtain the claimed geometric decrease

$$w_{t+1} \leq w_t \left(1 - \frac{\nu^2 \mu_{\mathcal{L}}}{2C_{\mathcal{L}}} \right). \quad (131)$$

If $\frac{\nu}{2C_{\mathcal{L}}}g_t > \gamma_{\max}$, then we have $\gamma_t = \gamma_{\max}$ and so (129) becomes:

$$\begin{aligned} w_{t+1} &\leq w_t - \nu\gamma_{\max}g_t + \gamma_{\max}^2 C_{\mathcal{L}} \\ &\leq w_t - \nu\gamma_{\max}g_t + \frac{\nu}{2}\gamma_{\max}g_t && (\text{using } C_{\mathcal{L}} < \frac{\nu}{2\gamma_{\max}}g_t) \end{aligned} \quad (132)$$

$$\begin{aligned} &\leq w_t - \frac{\nu}{2}\gamma_{\max}g_t \\ &\leq w_t \left(1 - \frac{\nu}{2}\gamma_{\max} \right). && (w_t \leq g_t \text{ by Lemma 17}) \end{aligned} \quad (133)$$

If $\gamma_{\max} \geq 1$ (either we are taking a FW step or an away step with a big step size), then the geometric rate is at least $(1 - \frac{\nu}{2})$, which is a better rate than $\rho_{\mathcal{L}}$ since $\nu^2 \leq \nu$ as $\nu \leq 1$, and one can show that $\frac{\mu_{\mathcal{L}}}{C_{\mathcal{L}}} \leq 1$ always (see Remark 7 in Appendix D of (Lacoste-Julien and Jaggi, 2015) for case (P) and use

a similar argument for case (I)). Thus $\rho_{\mathcal{L}}$ is valid both when $\gamma_t < \gamma_{\max}$ or $\gamma_{\max} \geq 1$, as claimed in the theorem.

When $\gamma_t = \gamma_{\max} < 1$, we cannot guarantee sufficient progress as γ_{\max} could be arbitrarily small (this can only happen for an away step as $\gamma_{\max} = 1$ for a FW step). These are the problematic *drop steps*, but as explained in Appendix A with Equation (38), they cannot happen too often for SP-AFW.

Finally, to show that the suboptimality cannot increase during a drop step ($\gamma_t = \gamma_{\max}$), we point out that the function $\gamma \mapsto w_t - \gamma \nu^{\text{PFW}} g_t^{\text{PFW}} + \gamma^2 C_{\mathcal{L}}$ is a convex function that is minimized by $\bar{\gamma} = \frac{\nu^{\text{PFW}}}{2C_{\mathcal{L}}} g_t^{\text{PFW}}$ and so is decreasing on $[0, \bar{\gamma}]$. When $\gamma_t = \gamma_{\max}$, we have that $\gamma_{\max} \leq \bar{\gamma}$, and thus the value for $\gamma = \gamma_{\max}$ is lower than the value for $\gamma = 0$, i.e.

$$w_{t+1} \leq w_t - \gamma_{\max} \nu^{\text{PFW}} g_t^{\text{PFW}} + \gamma_{\max}^2 C_{\mathcal{L}} \leq w_t. \quad (134)$$

□

The previous lemma (Lemma 23), the fact that the gap upper bounds the suboptimality (Lemma 17) and the primal suboptimality analysis lead us directly to the following theorem. This theorem is the affine invariant formulation with adaptive step size of Theorem 1.

Theorem 24. *Let \mathcal{L} be a strongly convex-concave function with a finite smoothness constant $C_{\mathcal{L}}$, a positive interior strong convex-concavity constant $\mu_{\mathcal{L}}^{\text{int}}$ (66) or a positive geometric strong convex-concavity $\mu_{\mathcal{L}}^{\text{A}}$ (71). Let us also define the rate multipliers ν as*

$$\nu^{\text{FW}} := 1 - \frac{M_{\mathcal{L}}}{\sqrt{\mu_{\mathcal{L}}^{\text{int}}}} \quad \text{and} \quad \nu^{\text{PFW}} := \frac{1}{2} - \frac{M_{\mathcal{L}}}{\sqrt{\mu_{\mathcal{L}}^{\text{A}}}} \quad (\text{see Equation (75) for the definition of } M_{\mathcal{L}}). \quad (135)$$

Let the tuple $(g_t, \nu, \mu_{\mathcal{L}})$ refers to either $(g_t^{\text{FW}}, \nu^{\text{FW}}, \mu_{\mathcal{L}}^{\text{int}})$ for case (I) where the algorithm is SP-FW, or $(g_t^{\text{PFW}}, \nu^{\text{PFW}}, \mu_{\mathcal{L}}^{\text{A}})$ for case (P) where the algorithm is SP-AFW,

If $\nu > 0$, then the suboptimality h_t of the iterates of the algorithm with step size $\gamma_t = \min(\gamma_{\max}, \frac{\nu}{2C_{\mathcal{L}}} g_t)$ decreases geometrically⁴ as,

$$h_t \leq P_{\mathcal{L}} \sqrt{2w_0} (1 - \rho_{\mathcal{L}})^{k(t)/2} \quad (136)$$

where $\rho_{\mathcal{L}} := \frac{\nu^2 \mu_{\mathcal{L}}}{2C_{\mathcal{L}}}$ and $k(t)$ is the number of non-drop step after t steps. For SP-FW, $k(t) = t$ and for SP-AFW, $k(t) \geq t/3$. Moreover we can also upper bound the minimum gap observed, for all $T \in \mathbb{N}$

$$\min_{t \leq T} g_t \leq \frac{2 \max\{\sqrt{C_{\mathcal{L}}}, \sqrt{w_0} (1 - \rho_{\mathcal{L}})^{k(T)/2}\}}{\nu} \sqrt{w_0} (1 - \rho_{\mathcal{L}})^{k(T)/2}. \quad (137)$$

The Theorem 1 statement can be deduced from this theorem using the lower and upper bounds on the affine invariant constant of this statement. More precisely, one can upper bound $C_{\mathcal{L}}$, $M_{\mathcal{L}}$, $P_{\mathcal{L}}$ respectively with Propositions 9, 13 and 14 and lower bound $\mu_{\mathcal{L}}^{\text{int}}$ and $\mu_{\mathcal{L}}^{\text{A}}$ respectively with Proposition 12 and Equation (72).⁵ If we apply these bounds to the rate multipliers in (135), it gives the bigger rate multipliers ν stated in Theorem 1.

Proof. We uses the Lemma 23 giving a geometric scheme, with a straightforward recurrence we prove that,

$$w_t \leq w_0 (1 - \rho_{\mathcal{L}})^{k(t)}, \quad (138)$$

⁴For a non-drop step one can use Corollary 22 to get the better rate on h_t losing the square root but with a potentially worse constant $h_t \leq \frac{2 \max\{C_{\mathcal{L}}, w_0 (1 - \rho_{\mathcal{L}})^{k(t)}\}}{\mu_{\mathcal{L}}^{\text{A}} (\nu^{\text{PFW}})^2} w_0 (1 - \rho_{\mathcal{L}})^{k(t)}$.

⁵Note that only the definition of δ_{μ} in Theorem 1 for case (P) requires to use the Euclidean norm (because inequality (72) with the pyramidal width only holds for the Euclidean norm). On the other hand, any norm could be used for (separately) bounding $C_{\mathcal{L}}$, $M_{\mathcal{L}}$ and $P_{\mathcal{L}}$.

where $k(t)$ is the number of non-drop step steps. This number is equal to t for the SP-FW algorithm and it is lower bounded by $t/3$ for the SP-AFW algorithm (see Section A Equation (38)). Then by using Proposition (15) relating h_t and the square root of w_t we get the first statement of the theorem,

$$h_t \leq P_{\mathcal{L}} \sqrt{2w_0} (1 - \rho_{\mathcal{L}})^{k(t)/2}. \quad (139)$$

To prove the second statement of the theorem we just use Theorem 21 for the last *non-drop step* after T iterations (let us assume it was at step t_0),

$$\begin{aligned} g_{t_0} &\leq \frac{2 \max\{\sqrt{C_{\mathcal{L}}}, \sqrt{w_0}(1 - \rho_{\mathcal{L}})^{k(t_0)/2}\}}{\nu} \sqrt{w_0}(1 - \rho_{\mathcal{L}})^{k(t_0)/2} \\ &= \frac{2 \max\{\sqrt{C_{\mathcal{L}}}, \sqrt{w_0}(1 - \rho_{\mathcal{L}})^{k(T)/2}\}}{\nu} \sqrt{w_0}(1 - \rho_{\mathcal{L}})^{k(T)/2}. \end{aligned} \quad (140)$$

(because $k(t_0) = k(T)$) (141)

The minimum of the gaps observed is smaller than the gap at time t_0 then,

$$\min_{t \leq T} g_t \leq g_{t_0} \leq \frac{2 \max\{\sqrt{C_{\mathcal{L}}}, \sqrt{w_0}(1 - \rho_{\mathcal{L}})^{k(T)/2}\}}{\nu} \sqrt{w_0}(1 - \rho_{\mathcal{L}})^{k(T)/2}. \quad (142)$$

□

The affine invariant formulation with the universal step size $\gamma_t = \min\left\{\gamma_{\max}, \frac{2}{2+k(t)}\right\}$ of Theorem 1 also follows from Lemma 23 by re-using standard FW proof patterns.

Theorem 25. *Let \mathcal{L} be a strongly convex-concave function with a finite smoothness constant $C_{\mathcal{L}}$, a positive interior strong convex-concavity constant $\mu_{\mathcal{L}}^{\text{int}}$ (66) or a positive geometric strong convex-concavity $\mu_{\mathcal{L}}^{\text{A}}$ (71). Let us also define the rate multipliers ν as*

$$\nu^{\text{FW}} := 1 - \frac{M_{\mathcal{L}}}{\sqrt{\mu_{\mathcal{L}}^{\text{int}}}} \quad \text{and} \quad \tilde{\nu}^{\text{PFW}} := 1 - \frac{M_{\mathcal{L}}}{\sqrt{\mu_{\mathcal{L}}^{\text{A}}}} \quad (\text{see Equation (75) for the definition of } M_{\mathcal{L}}). \quad (143)$$

Let ν refers to either ν^{FW} for case (I) where the algorithm is SP-FW, or $\tilde{\nu}^{\text{PFW}}$ for case (P) where the algorithm is SP-AFW,

If $\nu > \frac{1}{2}$, then the suboptimality w_t of the iterates of the algorithm with universal step size $\gamma_t = \min\left\{\gamma_{\max}, \frac{2}{2+k(t)}\right\}$ (see Equation (34) for more details about γ_{\max}) has the following decreasing upper bound:

$$w_t \leq \frac{C}{2 + k(t)} \quad (144)$$

where $C = 2 \max\left(w_0, \frac{2C_{\mathcal{L}}}{2\nu-1}\right)$ and $k(t)$ is the number of non-drop step after t steps. For SP-FW, $k(t) = t$ and for SP-AFW, $k(t) \geq t/3$. Moreover we can also upper bound the minimum FW gap observed for $T \geq 1$,

$$\min_{t \leq T} g_t^{\text{FW}} \leq \frac{5C}{\nu(k(T) + 1)}. \quad (145)$$

Note that in this theorem the constant $\tilde{\nu}^{\text{PFW}}$ is slightly different from the constant ν^{PFW} in Theorem 24.

Proof. We can put both the recurrence (115) for the SP-FW algorithm and the recurrence (117) for the SP-AFW algorithm (from the proof of Lemma 20) in the following form by using our unified notation introduced in the theorem statement:

$$w_{t+1} \leq w_t - \gamma_t \nu g_t + \gamma_t^2 C_{\mathcal{L}}. \quad (146)$$

Note that the gap g_t is the one defined in Equation (88) and depends on the algorithm. Let (ν) to stand respectively for (ν^{FW}) for SP-FW (case (I)) or $(\tilde{\nu}^{\text{PFW}})$ for SP-AFW (case (P)). With this notation, the inequality $g_t \geq w_t$ leads to,

$$w_{t+1} \leq w_t(1 - \nu\gamma_t) + \gamma_t^2 C_{\mathcal{L}}. \quad (147)$$

Our goal is to show by induction that

$$w_t \leq \frac{C}{2 + k(t)} \quad \text{where } C := 2 \max\left(w_0, \frac{2C_{\mathcal{L}}}{2\nu - 1}\right). \quad (\star)$$

Let us first define the convex function $f_t : \gamma \mapsto w_t(1 - \nu\gamma) + \gamma^2 C_{\mathcal{L}}$. We will show that under (\star) , the function f_t has the following property:

$$f_t\left(\frac{2}{2 + k(t)}\right) \leq \frac{C}{3 + k(t)}. \quad (148)$$

This property is due to a simple inequality on integers; let $k = k(t)$, from the crucial induction assumption, we get:

$$f_t\left(\frac{2}{2 + k}\right) = w_t \frac{2 + k - 2\nu}{2 + k} + \frac{4}{(2 + k)^2} C_{\mathcal{L}} \leq \frac{C}{3 + k} \left[\frac{(3 + k)(k + 1 - (2\nu - 1) + \frac{4C_{\mathcal{L}}}{C})}{(2 + k)^2} \right], \quad (149)$$

but $(2\nu - 1) \geq \frac{4C_{\mathcal{L}}}{C}$ and $(3 + k)(1 + k) < (2 + k)(2 + k)$ for any k , thus

$$f_t\left(\frac{2}{2 + k}\right) \leq \frac{C}{3 + k}. \quad (150)$$

Equation (150) is crucial for the inductive step of our recurrence.

- Hypothesis (\star) is true for $t = 0$ because $k(0) = 0$.
- Now let us assume that (\star) is true for a $t \in \mathbb{N}$. We set the stepsize $\gamma_t := \min\left\{\gamma_{\max}, \frac{2}{2 + k(t)}\right\}$.

If $k(t + 1) = k(t) + 1$, it means that $\gamma_t = \frac{2}{2 + k(t)}$ and then by (147) and (150),

$$w_{t+1} \leq f_t\left(\frac{2}{2 + k(t)}\right) \leq \frac{C}{3 + k(t)} = \frac{C}{2 + k(t + 1)}. \quad (151)$$

If $k(t + 1) = k(t)$, then it means that $0 \leq \gamma_t < \frac{2}{2 + k(t)}$. Hence, the convexity of the function f_t leads us to the inequality

$$\begin{aligned} w_{t+1} \leq f_t(\gamma_t) &\leq \max\left\{f_t(0), f_t\left(\frac{2}{2 + k(t)}\right)\right\} \\ &= \max\left\{w_t, f_t\left(\frac{2}{2 + k(t)}\right)\right\} \end{aligned} \quad (152)$$

$$\leq \max\left\{\frac{C}{2 + k(t)}, \frac{C}{3 + k(t)}\right\} \quad (153)$$

$$\leq \frac{C}{2 + k(t)}. \quad (154)$$

where we used (150) and the induction hypothesis (\star) to get the penultimate inequality (152). Since we assumed that $k(t + 1) = k(t)$, we get

$$w_{t+1} \leq \frac{C}{2 + k(t + 1)}, \quad (155)$$

completing the induction proof for (144).

In case (I), $k(t) = t$ and in case (P), $k(t) \geq t/3$ (see Equation (38)), leading us to the first statement of our theorem.

The proof of the second statement is inspired by the proof of Theorem C.3 from (Lacoste-Julien et al., 2013).

With the same notation as the proof of Lemma 20, we start from Equation (146) where we isolated the gap g_t to get the crucial inequality

$$g_t \leq \frac{w_t - w_{t+1}}{\nu\gamma_t} + \gamma_t \frac{C_{\mathcal{L}}}{\nu}. \quad (156)$$

Since the gap g_t is the one depending on the algorithm defined by $g_t := \langle -\mathbf{r}^{(t)}, \mathbf{d}^{(t)} \rangle$, we have $g_t = g_t^{\text{FW}}$ for SP-FW and $g_t = \max(g_t^{\text{FW}}, g_t^{\text{A}}) \geq g_t^{\text{FW}}$ for SP-AFW. Thus,

$$g_t^{\text{FW}} \leq g_t \leq \frac{w_t - w_{t+1}}{\nu\gamma_t} + \gamma_t \frac{C_{\mathcal{L}}}{\nu}. \quad (157)$$

In the following in order not to be too heavy with notation we will work with de FW gap and note g_t for g_t^{FW} .

The proof idea is to take a convex combination of the inequality (157) to obtain a new upper-bound on a convex combination of the gaps computed from step 0 to step T . Let us introduce the convex combination weight $\rho_t := \frac{\gamma_t \cdot k(t)(k(t)+2)}{S_T}$ where $k(t)$ is the number of non-drop steps after t steps and S_T is the normalization factor. Let us also call $N_T := \{t \leq T \mid t \text{ is a non-drop step}\}$. Taking the convex combination of (157), we get

$$\sum_{t=0}^T \rho_t g_t \leq \sum_{t=0}^T \rho_t \frac{w_t - w_{t+1}}{\nu\gamma_t} + \sum_{t=0}^T \rho_t \gamma_t \frac{C_{\mathcal{L}}}{\nu}. \quad (158)$$

By regrouping the terms and ignoring the negative term, we get

$$\sum_{t=0}^T \rho_t g_t \leq \frac{w_0 \rho_0}{\nu\gamma_0} + \frac{1}{\nu} \sum_{t=0}^{T-1} w_{t+1} \left(\frac{\rho_{t+1}}{\gamma_{t+1}} - \frac{\rho_t}{\gamma_t} \right) + \sum_{t=0}^T \rho_t \gamma_t \frac{C_{\mathcal{L}}}{\nu}. \quad (159)$$

By definition $\frac{\rho_t}{\gamma_t} := \frac{k(t)(k(t)+2)}{S_T}$ and notice that $\rho_0 = 0$. We now consider two possibilities: if γ_t is a drop step, then $k(t+1) = k(t)$ and so

$$\frac{\rho_{t+1}}{\gamma_{t+1}} - \frac{\rho_t}{\gamma_t} = 0. \quad (160)$$

If γ_t is a non-drop step, then $k(t+1) = k(t) + 1$ and thus we have

$$\frac{\rho_{t+1}}{\gamma_{t+1}} - \frac{\rho_t}{\gamma_t} = \frac{(k(t)+1)(k(t)+3)}{S_T} - \frac{k(t)(k(t)+2)}{S_T} \quad (161)$$

$$= \frac{2k(t)+3}{S_T}. \quad (162)$$

As $\gamma_t \leq \frac{2}{k(t)+2}$, we also have $\rho_t \gamma_t \leq \frac{4k(t)}{S_T(k(t)+2)}$. The normalization factor S_T to define a convex combination is equal to

$$S_T := \sum_{u=0}^T \gamma_u \cdot k(u)(k(u)+2) \geq \sum_{\substack{u=0 \\ u \in N_T}}^T \frac{2}{2+k(u)} \cdot k(u)(k(u)+2) = \sum_{k=0}^{k(T)} 2k = k(T)(K(T)+1). \quad (163)$$

Plugging this and (162) in the inequality (159) with the rate (\star) shown by induction gives us,

$$\sum_{t=0}^T \rho_t g_t \leq 0 + \frac{1}{\nu} \sum_{\substack{t=0 \\ t \in N_T}}^{T-1} \frac{C}{2+k(t+1)} \frac{2k(t)+3}{k(T)(k(T)+1)} + \sum_{t=0}^T \frac{4k(t)}{k(T)(k(T)+1)(k(t)+2)} \frac{C_{\mathcal{L}}}{\nu} \quad (164)$$

$$\leq \frac{2C}{\nu} \frac{1}{k(T)(k(T)+1)} \left(\sum_{\substack{t=0 \\ t \in N_T}}^{T-1} 1 + \frac{2C_{\mathcal{L}}}{C} \sum_{t=1}^T 1 \right) \quad (165)$$

$$\leq \frac{2C}{\nu(k(T)+1)} \left(1 + \frac{\nu}{2} \frac{T}{k(T)} \right) \leq \frac{5C}{\nu(k(T)+1)}, \quad (166)$$

by using $k(T) \geq T/3$. Finally, the minimum of the gaps is always smaller than any convex combination, so we can conclude that (for $T \geq 1$):

$$\min_{0 \leq t \leq T} g_t \leq \frac{5C}{\nu(k(T)+1)}. \quad (167)$$

□

E Strongly convex sets

In this section, we are going to prove that the function $\mathbf{s}(\cdot)$ is Lipschitz continuous when the sets \mathcal{X} and \mathcal{Y} are strongly convex and when the norm of the two gradient components are uniformly lower bounded. We will also give the details of the convergence rate proof for the strongly convex sets situation. Our proof uses similar arguments as Dunn (1979, Theorem 3.4 and 3.6).

Theorem' 3. *Let \mathcal{X} and \mathcal{Y} be β -strongly convex sets. If $\min(\|\nabla_x L(\mathbf{z})\|_{\mathcal{X}^*}, \|\nabla_y L(\mathbf{z})\|_{\mathcal{Y}^*}) \geq \delta > 0$ for all $\mathbf{z} \in \mathcal{X} \times \mathcal{Y}$, then the oracle function $\mathbf{z} \mapsto \mathbf{s}(\mathbf{z}) := \arg \min_{\mathbf{s} \in \mathcal{X} \times \mathcal{Y}} \langle \mathbf{s}, F(\mathbf{z}) \rangle$ is well defined and is $\frac{4L}{\delta\beta}$ -Lipschitz continuous (using the norm $\|(\mathbf{x}, \mathbf{y})\|_{\mathcal{X} \times \mathcal{Y}} := \|\mathbf{x}\|_{\mathcal{X}} + \|\mathbf{y}\|_{\mathcal{Y}}$, where $F(\mathbf{z}) := (\nabla_x \mathcal{L}(\mathbf{z}), -\nabla_y \mathcal{L}(\mathbf{z}))$).*

Proof. First note that since the sets are strongly convex, the minimum is reached at a unique point. Then, we introduce the following lemma which can be used to show that each component of the gradient is Lipschitz continuous irrespective of the other set.

Lemma 26. *Let $F_x : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}^d$ be a L -Lipschitz continuous function (i.e. $\|F_x(\mathbf{z}) - F_x(\mathbf{z}')\|_{\mathcal{X}^*} \leq L\|\mathbf{z} - \mathbf{z}'\|_{\mathcal{X} \times \mathcal{A}}$) and \mathcal{X} a β -strongly convex set. If $\forall \mathbf{z} \in \mathcal{X} \times \mathcal{A}$, $\|F_x(\mathbf{z})\|_{\mathcal{X}^*} \geq \delta > 0$, then $\mathbf{s}_x : \mathbf{z} \mapsto \arg \min_{\mathbf{s} \in \mathcal{X}} \langle \mathbf{s}, F_x(\mathbf{z}) \rangle$ is $\frac{2L}{\delta\beta}$ -Lipschitz continuous.*

Proof. Let $\mathbf{z}, \mathbf{z}' \in \mathcal{X} \times \mathcal{A}$ and let $\bar{\mathbf{x}} = \frac{\mathbf{s}_x(\mathbf{z}) + \mathbf{s}_x(\mathbf{z}')}{2}$, then

$$\begin{aligned} \langle \mathbf{s}_x(\mathbf{z}) - \mathbf{s}_x(\mathbf{z}'), -F_x(\mathbf{z}) \rangle &= 2 \langle \mathbf{s}_x(\mathbf{z}) - \bar{\mathbf{x}}, -F_x(\mathbf{z}) \rangle \\ &\geq 2 \langle \mathbf{x} - \bar{\mathbf{x}}, -F_x(\mathbf{z}) \rangle. \quad \forall \mathbf{x} \in \mathcal{X} \quad (\text{by definition of } s_x) \end{aligned} \quad (168)$$

Now (168) holds for any $\mathbf{x} \in B_\beta(\frac{1}{2}, \mathbf{s}_x(\mathbf{z}), \mathbf{s}_x(\mathbf{z}'))$ as this set is included in \mathcal{X} by β -strong convexity of \mathcal{X} . Then since $\bar{\mathbf{x}}$ is the center of $B_\beta(\frac{1}{2}, \mathbf{s}_x(\mathbf{z}), \mathbf{s}_x(\mathbf{z}'))$, we can choose a \mathbf{x} in this ball such that $\mathbf{x} - \bar{\mathbf{x}}$ is in the direction which achieves the dual norm of $-F_x(\mathbf{z})$.⁶ More specifically, we have that:

$$\| -F_x(\mathbf{z}) \|_{\mathcal{X}^*} = \sup_{\|\mathbf{v}\|_{\mathcal{X}} \leq 1} \langle -F_x(\mathbf{z}), \mathbf{v} \rangle.$$

⁶For the Euclidean norm, we choose $\mathbf{x} - \bar{\mathbf{x}}$ proportional to $-F_x$; but for general norms, it could be a different direction.

As we are in finite dimensions, this supremum is achieved by some vector \mathbf{v} . So choose $\mathbf{x} := \bar{\mathbf{x}} + \frac{\mathbf{v}}{\|\mathbf{v}\|_{\mathcal{X}}} \frac{\beta}{8} \|\mathbf{s}_x(\mathbf{z}) - \mathbf{s}_x(\mathbf{z}')\|_{\mathcal{X}}^2 \in B_\beta(\frac{1}{2}, \mathbf{s}_x(\mathbf{z}), \mathbf{s}_x(\mathbf{z}'))$ and plug it in (168):

$$\langle \mathbf{s}_x(\mathbf{z}) - \mathbf{s}_x(\mathbf{z}'), -F_x(\mathbf{z}) \rangle \geq \frac{\beta \|\mathbf{s}_x(\mathbf{z}) - \mathbf{s}_x(\mathbf{z}')\|_{\mathcal{X}}^2}{4 \|\mathbf{v}\|_{\mathcal{X}}} \langle \mathbf{v}, -F_x(\mathbf{z}) \rangle \quad (169)$$

$$= \frac{\beta \|\mathbf{s}_x(\mathbf{z}) - \mathbf{s}_x(\mathbf{z}')\|_{\mathcal{X}}^2}{4 \|\mathbf{v}\|_{\mathcal{X}}} \|F_x(\mathbf{z})\|_{\mathcal{X}^*} \quad (170)$$

$$\geq \frac{\beta}{4} \|\mathbf{s}_x(\mathbf{z}) - \mathbf{s}_x(\mathbf{z}')\|_{\mathcal{X}}^2 \|F_x(\mathbf{z})\|_{\mathcal{X}^*}. \quad (171)$$

Switching \mathbf{z} and \mathbf{z}' and using a similar argument, we get,

$$\langle \mathbf{s}_x(\mathbf{z}') - \mathbf{s}_x(\mathbf{z}), -F_x(\mathbf{z}') \rangle \geq \frac{\beta}{4} \|\mathbf{s}_x(\mathbf{z}) - \mathbf{s}_x(\mathbf{z}')\|_{\mathcal{X}}^2 \|F_x(\mathbf{z}')\|_{\mathcal{X}^*}. \quad (172)$$

Hence summing (171) and (172),

$$\begin{aligned} \frac{\beta}{4} (\|F_x(\mathbf{z})\|_{\mathcal{X}^*} + \|F_x(\mathbf{z}')\|_{\mathcal{X}^*}) \|\mathbf{s}_x(\mathbf{z}) - \mathbf{s}_x(\mathbf{z}')\|_{\mathcal{X}}^2 &\leq \langle \mathbf{s}_x(\mathbf{z}) - \mathbf{s}_x(\mathbf{z}'), F_x(\mathbf{z}') - F_x(\mathbf{z}) \rangle \\ &\leq \|\mathbf{s}_x(\mathbf{z}) - \mathbf{s}_x(\mathbf{z}')\|_{\mathcal{X}} \|F_x(\mathbf{z}') - F_x(\mathbf{z})\|_{\mathcal{X}^*} \\ &\leq L \|\mathbf{s}_x(\mathbf{z}) - \mathbf{s}_x(\mathbf{z}')\|_{\mathcal{X}} \|\mathbf{z}' - \mathbf{z}\|_{\mathcal{X} \times \mathcal{Y}} \quad (\text{Lip. cty. of } F_x) \end{aligned} \quad (173)$$

and finally

$$\|\mathbf{s}_x(\mathbf{z}) - \mathbf{s}_x(\mathbf{z}')\|_{\mathcal{X}} \leq \frac{4L}{\beta (\|F_x(\mathbf{z})\|_{\mathcal{X}^*} + \|F_x(\mathbf{z}')\|_{\mathcal{X}^*})} \|\mathbf{z} - \mathbf{z}'\|_{\mathcal{X} \times \mathcal{Y}} \leq \frac{2L}{\delta\beta} \|\mathbf{z} - \mathbf{z}'\|_{\mathcal{X} \times \mathcal{Y}}. \quad (174)$$

□

To prove our theorem, we will notice that for the saddle point setup, the oracle function $\mathbf{s}(\cdot) := \arg \min_{\mathbf{s} \in \mathcal{X} \times \mathcal{Y}} \langle \mathbf{s}, F(\cdot) \rangle$ can be decomposed as $\mathbf{s}(\cdot) = (\mathbf{s}_x(\cdot), \mathbf{s}_y(\cdot))$ where $\mathbf{s}_x(\cdot) := \arg \min_{\mathbf{s} \in \mathcal{X}} \langle \mathbf{s}, F_x(\cdot) \rangle$ and $\mathbf{s}_y(\cdot) := \arg \min_{\mathbf{s} \in \mathcal{Y}} \langle \mathbf{s}, F_y(\cdot) \rangle$. Then applying our lemma, the function $\mathbf{s}_x(\cdot)$ is Lipschitz continuous. The same way $\mathbf{s}_y(\cdot)$ is Lipschitz continuous. Then, for all \mathbf{z}, \mathbf{z}' in $\mathcal{X} \times \mathcal{Y}$

$$\|\mathbf{s}(\mathbf{z}) - \mathbf{s}(\mathbf{z}')\|_{\mathcal{X} \times \mathcal{Y}} = \|\mathbf{s}_x(\mathbf{z}) - \mathbf{s}_x(\mathbf{z}')\|_{\mathcal{X}} + \|\mathbf{s}_y(\mathbf{z}) - \mathbf{s}_y(\mathbf{z}')\|_{\mathcal{Y}} \leq \frac{4L}{\delta\beta} \|\mathbf{z} - \mathbf{z}'\|_{\mathcal{X} \times \mathcal{Y}}, \quad (175)$$

which gives the definition of the Lipschitz continuity of our function and proves the theorem. □

In this theorem, we introduced the function F . This function is monotone in the following sense:

$$\forall \mathbf{z}, \mathbf{z}' \quad \langle \mathbf{z} - \mathbf{z}', F(\mathbf{z}) - F(\mathbf{z}') \rangle \geq 0. \quad (176)$$

Actually this property follows directly from the convexity of $\mathcal{L}(\cdot, \mathbf{y})$ and the concavity of $\mathcal{L}(\mathbf{x}, \cdot)$. We can also prove that when the sets \mathcal{X} and \mathcal{Y} are strongly convex and when the gradient is uniformly lower bounded, we can relate the gap and the distance between $\mathbf{z}^{(t)}$ and $\mathbf{s}^{(t)}$.

Lemma 27. *If \mathcal{X} is a β -strongly convex set and if $\|\nabla f\|_{\mathcal{X}^*}$ is uniformly lower bounded by δ on \mathcal{X} , then*

$$\max_{\mathbf{s} \in \mathcal{X}} \langle \mathbf{s} - \mathbf{x}, -\nabla f(\mathbf{x}) \rangle \geq \frac{\beta}{4} \delta \|\mathbf{s}(\mathbf{x}) - \mathbf{x}\|^2, \quad (177)$$

where $\mathbf{s}(\mathbf{x}) = \arg \max_{\mathbf{s} \in \mathcal{X}} \langle \mathbf{s} - \mathbf{x}, -\nabla f(\mathbf{x}) \rangle$.

Proof. Let \mathbf{x} and $\mathbf{s}(\mathbf{x})$ be in \mathcal{X} . We have $B_\beta(\frac{1}{2}, \mathbf{s}(\mathbf{x}), \mathbf{x}) \subset \mathcal{X}$ by β -strong convexity. So as in the proof of Lemma 26, let \mathbf{v} be the vector such that $\|\mathbf{v}\|_{\mathcal{X}} \leq 1$ and $\langle -\nabla f(\mathbf{x}), \mathbf{v} \rangle = \|\nabla f(\mathbf{x})\|_{\mathcal{X}^*}$. Let

$$\bar{\mathbf{s}} := \frac{\mathbf{s}(\mathbf{x}) + \mathbf{x}}{2} + \frac{\beta}{8} \|\mathbf{s}(\mathbf{x}) - \mathbf{x}\|^2 \frac{\mathbf{v}}{\|\mathbf{v}\|_{\mathcal{X}}} \in \mathcal{X}. \quad (178)$$

Then

$$\begin{aligned} \langle \mathbf{s}(\mathbf{x}) - \mathbf{x}, -\nabla f(\mathbf{x}) \rangle &\geq \langle \bar{\mathbf{s}} - \mathbf{x}, -\nabla f(\mathbf{x}) \rangle \\ &= \frac{1}{2} \langle \mathbf{s}(\mathbf{x}) - \mathbf{x}, -\nabla f(\mathbf{x}) \rangle + \frac{\beta}{8} \|\mathbf{s}(\mathbf{x}) - \mathbf{x}\|^2 \frac{\|\nabla f(\mathbf{x})\|_{\mathcal{X}^*}}{\|\mathbf{v}\|_{\mathcal{X}}} \\ &\geq \frac{1}{2} \langle \mathbf{s}(\mathbf{x}) - \mathbf{x}, -\nabla f(\mathbf{x}) \rangle + \frac{\beta}{8} \delta \|\mathbf{s}(\mathbf{x}) - \mathbf{x}\|^2 \end{aligned} \quad (179)$$

which leads us to the desired result. \square

From this lemma, under the assumption that $\min(\|\nabla_x L(\mathbf{z})\|_{\mathcal{X}^*}, \|\nabla_y L(\mathbf{z})\|_{\mathcal{Y}^*}) \geq \delta \forall \mathbf{z} \in \mathcal{X} \times \mathcal{Y}$, it directly follows that

$$\begin{aligned} g_t^{\text{FW}} = g_t^{(x)} + g_t^{(y)} &\geq \frac{\beta}{4} \delta \left(\|\mathbf{s}_x^{(t)} - \mathbf{x}^{(t)}\|_{\mathcal{X}}^2 + \|\mathbf{s}_y^{(t)} - \mathbf{y}^{(t)}\|_{\mathcal{Y}}^2 \right) \\ &\geq \frac{\beta}{8} \delta \left(\|\mathbf{s}_x^{(t)} - \mathbf{x}^{(t)}\|_{\mathcal{X}} + \|\mathbf{s}_y^{(t)} - \mathbf{y}^{(t)}\|_{\mathcal{Y}} \right)^2 = \frac{\beta}{8} \delta \|\mathbf{s}^{(t)} - \mathbf{z}^{(t)}\|_{\mathcal{X} \times \mathcal{Y}}^2. \end{aligned} \quad (180)$$

Now we recall the convergence theorem for strongly convex sets from the main text, Theorem 4:

Theorem' 4. *Let \mathcal{L} be a convex-concave function and \mathcal{X} and \mathcal{Y} two compact β -strongly convex sets. Assume that the gradient of \mathcal{L} is L -Lipschitz continuous and that there exists $\delta > 0$ such that $\min(\|\nabla_x \mathcal{L}(\mathbf{z})\|_{\mathcal{X}^*}, \|\nabla_y \mathcal{L}(\mathbf{z})\|_{\mathcal{Y}^*}) \geq \delta \forall \mathbf{z} \in \mathcal{X} \times \mathcal{Y}$. Set $C_\delta := 2L + \frac{8L^2}{\beta\delta}$. Then the gap g_t^{FW} (7) of the SP-FW algorithm with step size $\gamma_t = \frac{g_t^{\text{FW}}}{\|\mathbf{s}^{(t)} - \mathbf{z}^{(t)}\|_{\mathcal{X} \times \mathcal{Y}}^2 C_\delta}$ converges linearly as*

$$g_t^{\text{FW}} \leq g_0 (1 - \rho)^t \quad (181)$$

where $\rho := \frac{\beta\delta}{16C_\delta}$. The initial gap g_0 is cheaply computed during the first step of the SP-FW algorithm. Alternatively, one can use the following upper bound to get uniform guarantees:

$$g_0 \leq \sup_{\mathbf{z} \in \mathcal{X} \times \mathcal{Y}} \|\nabla_x \mathcal{L}(\mathbf{z})\|_{\mathcal{X}^*} D_{\mathcal{X}} + \sup_{\mathbf{z} \in \mathcal{X} \times \mathcal{Y}} \|\nabla_y \mathcal{L}(\mathbf{z})\|_{\mathcal{Y}^*} D_{\mathcal{Y}}. \quad (182)$$

Proof. We compute the following relation on the gap:

$$\begin{aligned} g_{t+1} &= \left\langle \mathbf{z}^{(t+1)} - \mathbf{s}^{(t+1)}, F(\mathbf{z}^{(t+1)}) \right\rangle \\ &= \left\langle \mathbf{z}^{(t)} - \mathbf{s}^{(t+1)}, F(\mathbf{z}^{(t+1)}) \right\rangle + \gamma_t \left\langle \mathbf{s}^{(t)} - \mathbf{z}^{(t)}, F(\mathbf{z}^{(t+1)}) \right\rangle \\ &= \left\langle \mathbf{z}^{(t)} - \mathbf{s}^{(t+1)}, F(\mathbf{z}^{(t)}) \right\rangle + \left\langle \mathbf{z}^{(t)} - \mathbf{s}^{(t+1)}, F(\mathbf{z}^{(t+1)}) - F(\mathbf{z}^{(t)}) \right\rangle \\ &\quad + \gamma_t \left\langle \mathbf{s}^{(t)} - \mathbf{z}^{(t)}, F(\mathbf{z}^{(t)}) \right\rangle + \gamma_t \left\langle \mathbf{s}^{(t)} - \mathbf{z}^{(t)}, F(\mathbf{z}^{(t+1)}) - F(\mathbf{z}^{(t)}) \right\rangle \\ &\leq \left\langle \mathbf{z}^{(t)} - \mathbf{s}^{(t+1)}, F(\mathbf{z}^{(t)}) \right\rangle + \left\langle \mathbf{z}^{(t)} - \mathbf{s}^{(t+1)}, F(\mathbf{z}^{(t+1)}) - F(\mathbf{z}^{(t)}) \right\rangle \\ &\quad + \gamma_t \left\langle \mathbf{s}^{(t)} - \mathbf{z}^{(t)}, F(\mathbf{z}^{(t)}) \right\rangle + \gamma_t^2 \|\mathbf{s}^{(t)} - \mathbf{z}^{(t)}\|^2 L \end{aligned} \quad (183)$$

where in the last line we used the fact that the function $F(\cdot)$ is Lipschitz continuous. Then using that $\langle \mathbf{z}^{(t)} - \mathbf{s}^{(t+1)}, F(\mathbf{z}^{(t)}) \rangle \leq \langle \mathbf{z}^{(t)} - \mathbf{s}^{(t)}, F(\mathbf{z}^{(t)}) \rangle$ (by definition of $\mathbf{s}^{(t)}$), we get

$$\begin{aligned} g_{t+1} &\leq g_t(1 - \gamma_t) + \left\langle \mathbf{z}^{(t)} - \mathbf{s}^{(t+1)}, F(\mathbf{z}^{(t+1)}) - F(\mathbf{z}^{(t)}) \right\rangle + \gamma_t^2 \|\mathbf{s}^{(t)} - \mathbf{z}^{(t)}\|^2 L \\ &\leq g_t(1 - \gamma_t) + \left\langle \mathbf{s}^{(t)} - \mathbf{s}^{(t+1)}, F(\mathbf{z}^{(t+1)}) - F(\mathbf{z}^{(t)}) \right\rangle + \gamma_t^2 \|\mathbf{s}^{(t)} - \mathbf{z}^{(t)}\|^2 L. \end{aligned} \quad (184)$$

The last line uses the fact that F is monotone by convexity (Equation (176)). Finally, using once again the Lipschitz continuity of F and the one of $\mathbf{s}(\cdot)$ (by Theorem 3), we get

$$\begin{aligned} \left\langle \mathbf{s}^{(t)} - \mathbf{s}^{(t+1)}, F(\mathbf{z}^{(t+1)}) - F(\mathbf{z}^{(t)}) \right\rangle &\leq \|\mathbf{s}^{(t)} - \mathbf{s}^{(t+1)}\| L \|\mathbf{z}^{(t+1)} - \mathbf{z}^{(t)}\| \\ &\leq \frac{4L^2}{\beta\delta} \|\mathbf{z}^{(t+1)} - \mathbf{z}^{(t)}\|^2 \quad (\text{Lipschitz continuity of } \mathbf{s}) \\ &= \frac{4L^2}{\beta\delta} \gamma_t^2 \|\mathbf{s}^{(t)} - \mathbf{z}^{(t)}\|^2. \end{aligned} \quad (185)$$

Combining (185) with (184), we get

$$g_{t+1} \leq g_t(1 - \gamma_t) + \gamma_t^2 \|\mathbf{s}^{(t)} - \mathbf{z}^{(t)}\|^2 \frac{C_\delta}{2} \quad \text{where} \quad C_\delta := 2L + \frac{8L^2}{\beta\delta}. \quad (186)$$

Thus by setting the step size $\gamma_t = \frac{g_t}{\|\mathbf{s}^{(t)} - \mathbf{z}^{(t)}\|^2 C_\delta}$, we get

$$g_{t+1} \leq g_t - \frac{g_t}{2C_\delta} \left(\frac{g_t}{\|\mathbf{s}^{(t)} - \mathbf{z}^{(t)}\|^2} \right) \leq g_t \left(1 - \frac{\beta\delta}{16C_\delta} \right), \quad (187)$$

using the fact that the gap is lower bounded by a constant times the square of the distance between $\mathbf{s}^{(t)}$ and $\mathbf{z}^{(t)}$ (Equation (180)). \square

Note that the bound in this theorem is not affine invariant because of the presence of Lipschitz constants and strong convexity constants of the sets. The algorithm is not affine invariant either because the step size rule depends on these constants as well as on $\|\mathbf{s}^{(t)} - \mathbf{z}^{(t)}\|$. Deriving an affine invariant step size choice and convergence analysis is still an interesting open problem in this setting.

F Details on the experiments

Graphical Games. The payoff matrix M that we use encodes the following simple model of competition between universities with their respective benefits:

1. University 1 (respectively University 2) has benefit $b_i^{(1)}$ ($b_i^{(2)}$) to get student i .
2. Student i ranks the possible roommates with a permutation $\sigma_i \in \mathcal{S}_p$. Let $\sigma_i(j)$ represents the rank of j for i (first in the list is the preferred one).
3. They go to the university that matched them with their preferred roommate, in case of equality the student chooses randomly.
4. Supposing that \mathbf{x} encodes the roommate assignment proposed by University 1 (and \mathbf{y} for University 2), then the expectation of the benefit of University 1 is $\mathbf{x}^\top M \mathbf{y}$, with the following definition for the payoff matrix M indexed by pairs of matched students. For the pairs (i, j) with $i < j$ and (k, l) with $k < l$ with elements in $1, \dots, s$, we have:

$$\begin{aligned}
 \text{(a) } M_{ij,il} &= \begin{cases} b_i^{(1)} & \text{if } \sigma_i(j) < \sigma_i(l) \quad \text{i.e. student } i \text{ preferred } j \text{ over } l \\ -b_i^{(2)} & \text{if } \sigma_i(j) > \sigma_i(l) \\ \frac{b_i^{(1)} - b_i^{(2)}}{2} & \text{otherwise (in that case } j = l). \end{cases} \\
 \text{(b) } M_{ij,kj} &= M_{ji,jk} \\
 \text{(c) } M_{ij,ki} &= M_{ij,ik} \\
 \text{(d) } M_{ij,jl} &= M_{ij,lj} \\
 \text{(e) } M_{ij,kl} &= 0 \quad \text{otherwise}
 \end{aligned}$$

Note that we need to do unipartite matching here (and not bipartite matching) since we have to match students together and not students with dorms.

For our experiments, in order to get a realistic payoff matrix, we set $\mu_i \sim \mathcal{U}[0, 1]$ the *true* value of student i . Then we set $b_i^{(U)} \sim \mathcal{N}(\mu_i, 0.1)$ the value of the student i *observed* by University U . To solve the perfect matching problem, we used Blossom V by [Kolmogorov \(2009\)](#).

Sparse structured SVM. We give here more details on the derivations of the objective function for the structured SVM problem. We first recall the structured prediction setup with the same notation from [\(Lacoste-Julien et al., 2013\)](#). In structured prediction, the goal is to predict a structured object $\mathbf{y} \in \mathcal{Y}(\mathbf{x})$ (such as a sequence of tags) for a given input $\mathbf{x} \in \mathcal{X}$. For the structured SVM approach, a structured feature map $\phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$ encodes the relevant information for input / output pairs, and a linear classifier with parameter \mathbf{w} is defined by $h_{\mathbf{w}} = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \langle \mathbf{w}, \phi(\mathbf{x}, \mathbf{y}) \rangle$. We are also given a task-dependent structured error $L(\mathbf{y}', \mathbf{y})$ that gives the loss of predicting \mathbf{y} when the ground truth is \mathbf{y}' . Given a labeled training set $\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^n$, the standard ℓ_2 -regularized structured SVM objective in its non-smooth formulation for learning as given for example in Equation (3) from [\(Lacoste-Julien et al., 2013\)](#) is:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{n} \sum_i \tilde{H}_i(\mathbf{w}) \quad (188)$$

where $\tilde{H}_i(\mathbf{w}) := \max_{\mathbf{y} \in \mathcal{Y}_i} L_i(\mathbf{y}) - \langle \mathbf{w}, \psi_i(\mathbf{y}) \rangle$ is the structured hinge loss, and the following notational shorthands were defined: $\mathcal{Y}_i := \mathcal{Y}(\mathbf{x}^{(i)})$, $L_i(\mathbf{y}) := L(\mathbf{y}^{(i)}, \mathbf{y})$ and $\psi_i(\mathbf{y}) := \phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) - \phi(\mathbf{x}^{(i)}, \mathbf{y})$.

In our setting, we consider a sparsity inducing ℓ_1 -regularization instead. Moreover, we use the (equivalent) constrained formulation instead of the penalized one, in order to get a problem over a polytope. We thus get the following challenging problem:

$$\min_{\|\mathbf{w}\|_1 \leq R} \frac{1}{n} \sum_i \tilde{H}_i(\mathbf{w}). \quad (189)$$

To handle any type of structured output space \mathcal{Y} , we use the following generic encoding. Enumerating the elements of \mathcal{Y}_i , we can represent the j^{th} element of \mathcal{Y}_i as $(\overbrace{0, \dots, 0}^{j-1}, 1, 0, \dots, 0) \in \mathbb{R}^{|\mathcal{Y}_i|}$. Let M_i have $(\psi_i(\mathbf{y}))_{\mathbf{y} \in \mathcal{Y}_i}$ as columns and let \mathbf{L}_i be a vector of length $|\mathcal{Y}_i|$ with $L_i(\mathbf{y})$ as its entries. The functions $\tilde{H}_i(\mathbf{w})$ can then be rewritten as the maximization of linear functions in \mathbf{y} : $\tilde{H}_i(\mathbf{w}) = \max_{\mathbf{y} \in \mathcal{Y}_i} \mathbf{L}_i^\top \mathbf{y} - \mathbf{w}^\top M_i \mathbf{y}$. As the maximization of linear functions over a polytope is always obtained at one of its vertex, we can equivalently define the maximization over the convex hull of \mathcal{Y}_i , which is the probability simplex in $\mathbb{R}^{|\mathcal{Y}_i|}$ that we denote $\Delta(|\mathcal{Y}_i|)$:

$$\max_{\mathbf{y}_i \in \mathcal{Y}_i} \mathbf{L}_i^\top \mathbf{y}_i - \mathbf{w}^\top M_i \mathbf{y}_i = \max_{\alpha_i \in \Delta(|\mathcal{Y}_i|)} \mathbf{L}_i^\top \alpha_i - \mathbf{w}^\top M_i \alpha_i \quad (190)$$

Thus our equivalent objective is

$$\min_{\|\mathbf{w}\|_1 \leq R} \frac{1}{n} \sum_i \left(\max_{\mathbf{y}_i \in \mathcal{Y}_i} \mathbf{L}_i^\top \mathbf{y}_i - \mathbf{w}^\top M_i \mathbf{y}_i \right) = \min_{\|\mathbf{w}\|_1 \leq R} \frac{1}{n} \sum_i \left(\max_{\alpha_i \in \Delta(|\mathcal{Y}_i|)} \mathbf{L}_i^\top \alpha_i - \mathbf{w}^\top M_i \alpha_i \right), \quad (191)$$

which is the bilinear saddle point formulation given in the main text in (31).

Supplementary References

- V. Kolmogorov. Blossom V: a new implementation of a minimum cost perfect matching algorithm. *Mathematical Programming Computation*, 2009.
- S. Lacoste-Julien and M. Jaggi. An affine invariant linear convergence analysis for Frank-Wolfe algorithms. *arXiv preprint arXiv:1312.7864*, 2013.
- Y. Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer, 2004.
- A. Osokin, J.-B. Alayrac, I. Lukasewitz, P. K. Dokania, and S. Lacoste-Julien. Minding the gaps for block Frank-Wolfe optimization of structured SVMs. In *ICML*, 2016.