

Overview

Summary

Frank-Wolfe algorithm (FW) **gained in popularity** in the last couple of years because of some **key properties** allowing it to cheaply exploit the **structured constraint** sets appearing in machine learning applications.

- We tried to **extend** FW nice properties to solve saddle point problem.
- **Straightforward** extension but **Non trivial** analysis.

Contributions

- **Extend** several variants of the FW algorithm to solve SP problem.
- **Prove convergence** results for these methods over **polytope** domains giving a partial answer to **Hammond's conjecture** [1].
- **Introduce** application with fast linear minimization oracle (LMO)

Call for applications

Saddle Point (SP) Problem

Let $\mathcal{L} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, where \mathcal{X} and \mathcal{Y} are convex and compact.

Saddle point problem solve $\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} \mathcal{L}(x, y)$

A solution $z^* = (x^*, y^*)$ is called a **Saddle Point**.

Stationary conditions

$$\begin{cases} \langle x - x^*, \nabla_x \mathcal{L}(x^*, y^*) \rangle \geq 0 \\ \langle y - y^*, -\nabla_y \mathcal{L}(x^*, y^*) \rangle \geq 0 \end{cases}$$

Variational inequality

$$\begin{cases} \langle z - z^*, r(z^*) \rangle \geq 0, \forall z = (x, y) \\ r(z) = (\nabla_x \mathcal{L}(z), -\nabla_y \mathcal{L}(z)) \end{cases}$$

Global solution if \mathcal{L} convex-concave: $\forall (x, y) \in \mathcal{X} \times \mathcal{Y}$

$x' \mapsto \mathcal{L}(x', y)$ is convex and $y' \mapsto \mathcal{L}(x, y')$ is concave.

Standard method to solve SP

projected gradient

Simple algorithm to solve Saddle point optimization:

$$\begin{aligned} x^+ &= P_{\mathcal{X}}(x - \gamma \nabla_x \mathcal{L}(x, y)) \\ y^+ &= P_{\mathcal{Y}}(y + \gamma \nabla_y \mathcal{L}(x, y)) \end{aligned}$$

Non-smooth function:

$$\frac{1}{T} \sum_{t=1}^T (x^{(t)}, y^{(t)}) \xrightarrow{T \rightarrow \infty} (x^*, y^*)$$

projected extra-gradient

$$\begin{aligned} \bar{x} &= P_{\mathcal{X}}(x - \gamma \nabla_x \mathcal{L}(x, y)) \\ \bar{y} &= P_{\mathcal{Y}}(y + \gamma \nabla_y \mathcal{L}(x, y)) \\ x^+ &= P_{\mathcal{X}}(x - \gamma \nabla_x \mathcal{L}(\bar{x}, \bar{y})) \\ y^+ &= P_{\mathcal{Y}}(y + \gamma \nabla_y \mathcal{L}(\bar{x}, \bar{y})) \end{aligned}$$

Faster for Smooth function:

$$(x^{(T)}, y^{(T)}) \xrightarrow{T \rightarrow \infty} (x^*, y^*)$$

SP-FW

Update rule

Update

$$\begin{aligned} r(z) &:= \begin{pmatrix} \nabla_x \mathcal{L}(z) \\ -\nabla_y \mathcal{L}(z) \end{pmatrix} \\ s &\in \arg \min_{s' \in \mathcal{X} \times \mathcal{Y}} \langle s', r(z) \rangle \\ z^+ &:= (1 - \gamma)z + \gamma s \end{aligned}$$

Stopping criterion, the **FW gap** :

$$g_t := \langle r(z), z - s \rangle \leq \varepsilon$$

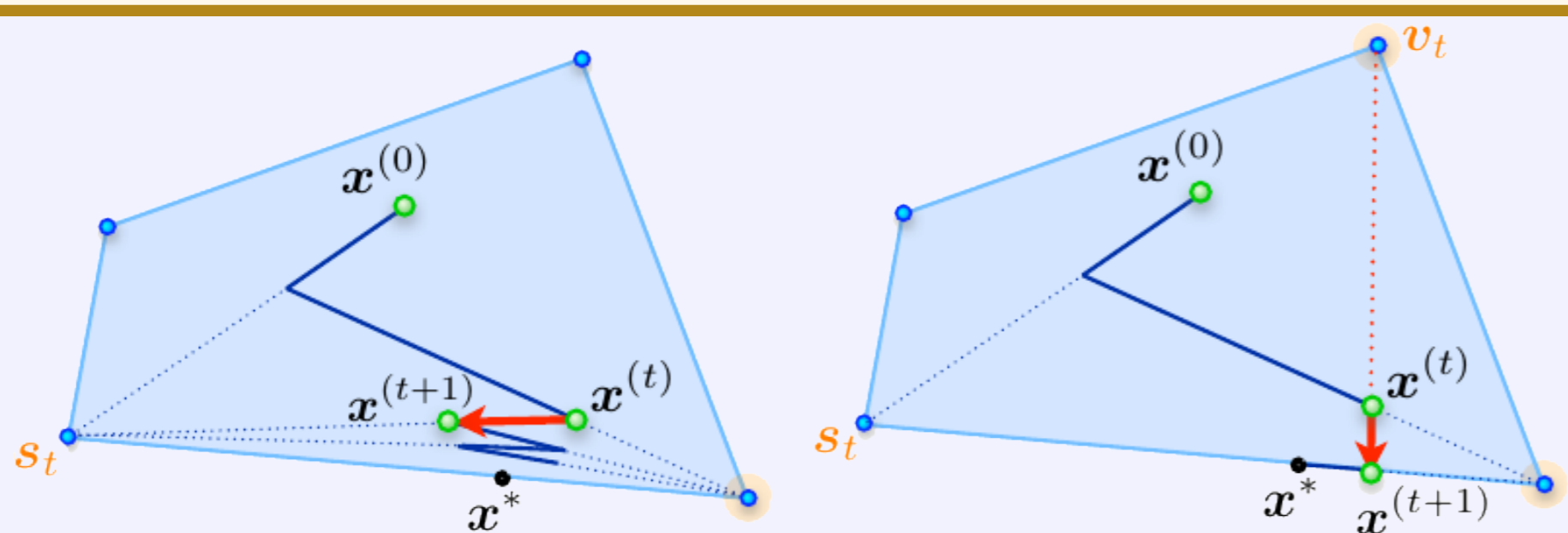
Properties

Only LMO.

- Gap certificate for free.
- Algorithm *affine invariant*.
- Universal step size $\gamma_t := \frac{2}{2+t}$ and adaptive step size $\gamma_t \sim g_t$.
- **No** line-search.
- **Sparsity** of the iterates.

Away step

$$v \in \arg \max_{v' \in \mathcal{S}_x \times \mathcal{S}_y} \langle r(z), v' \rangle, \quad d_A := z - v \text{ and } g_t^A := \langle d_A, -r(z) \rangle$$



Theoretical contribution

Convergence

SP extension of FW with **away step**:

- **Linear** rate with **adaptive** step size $\gamma_t := \frac{\nu}{LD^2} g_t$.

$$\min_{s \leq t} g_s \leq \left(1 - \nu^2 \frac{\delta^2 \mu}{D^2 2L}\right)^t$$

- **Sublinear** rate with **universal** step size $\gamma_t := \frac{2}{2+k(t)}$.

Hypothesis

Similar hypothesis as AFW:

- L -smooth, μ -strongly convex.
- \mathcal{X} and \mathcal{Y} **polytopes**.

Additional assumption on **bilinearity**:

$$\nu := \frac{1}{2} - \frac{\sqrt{2} \|M\| D}{\mu \delta} > 0$$

Details on the additional assumption

"Strong convexity μ big enough compared to the bilinear coupling $\|M\|$ "

$$\mathcal{L}(x, y) = f(x) + x^T M y - g(y).$$

$$D := \max\{\text{diam}(\mathcal{X}), \text{diam}(\mathcal{Y})\}, \quad \delta := \min\{PWidth(\mathcal{X}), PWidth(\mathcal{Y})\}$$

Difficulties for SP

FW proof technique

$$\begin{aligned} h_t &:= f(x^{(t)}) - \min_{x \in \mathcal{X}} f(x) \\ h_{t+1} &\leq h_t - \gamma_t g_t + \gamma_t^2 \frac{L \|d^{(t)}\|^2}{2} \end{aligned}$$

Set $\gamma_t = \frac{g_t}{C}$, decreasing scheme:

$$h_{t+1} \leq h_t - \frac{g_t^2}{2C} \leq h_t - \frac{h_t^2}{2C}$$

Same derivation for SP

$$\mathcal{L}_{t+1} \leq \mathcal{L}_t - \gamma_t \underbrace{(g_t^{(x)} - g_t^{(y)})}_{\text{arbitrary sign}} + \gamma_t^2 \frac{LD^2}{2}$$

- Cannot control oscillations of \mathcal{L}_t .
- Must introduce other quantities.
- Proof use **recent advances** on AFW [4].

Related conjectures

Karlin's conjecture

SP-FW is equivalent to the **fictitious play algorithm** [5] when

$$\gamma_t = \frac{1}{1+t} \text{ and } \mathcal{L}(x, y) = x^T M y$$

Karlin [3] conjectured that:

$$g_t \leq O\left(\frac{1}{\sqrt{t}}\right)$$

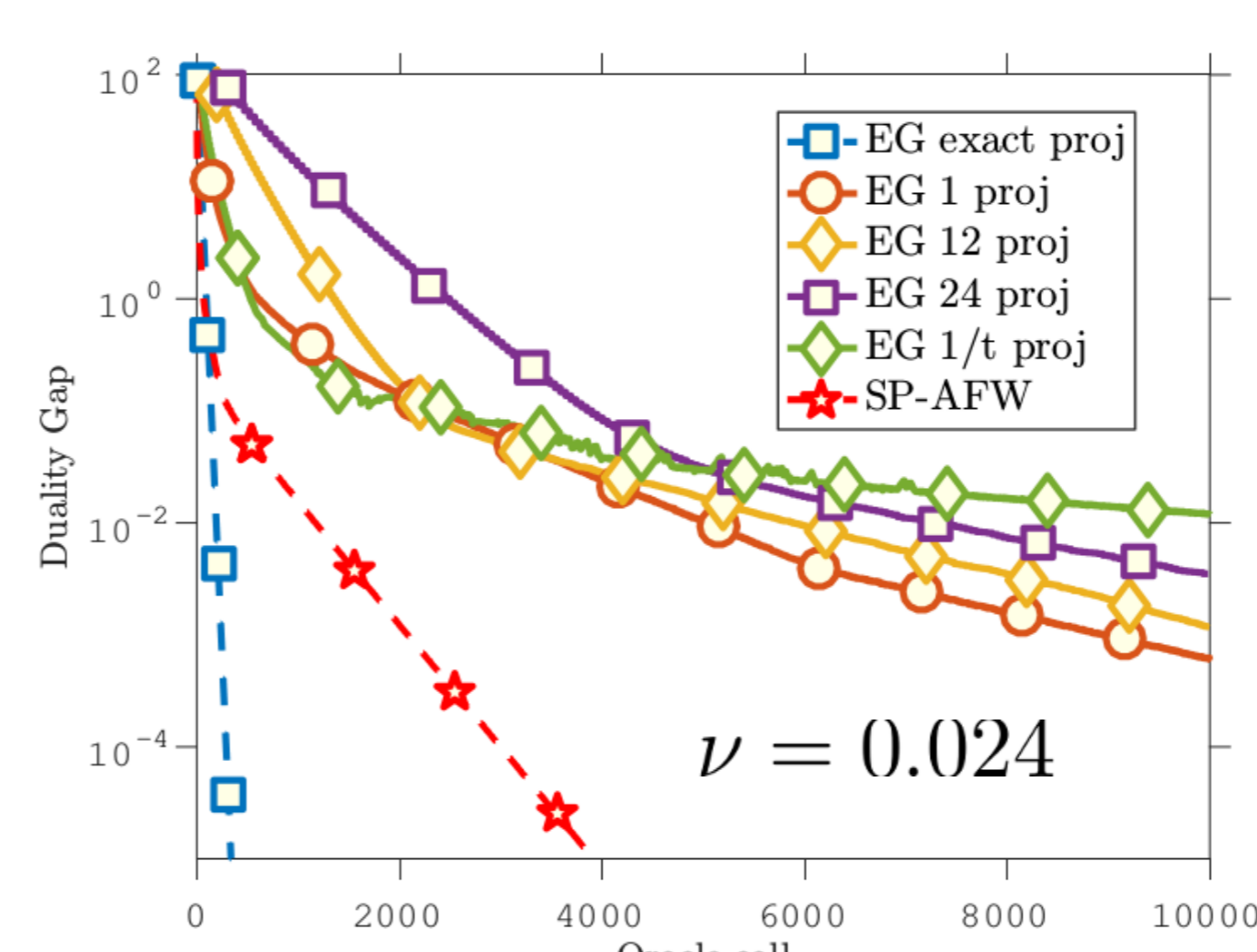
Hammond's conjecture

Hammond [1] conjectured that for Variational inequalities:

If g is **uniformly** monotone and the constraints is a **bounded polyhedron**, then the **fictitious play algorithm** will solve the variational inequality problem.

Toy experiments

For the Poster: we compared with EG by He and Harchaoui [2] (fig. on the left) performing projected extra-gradient with approximate projections.



Theoretical:

$$\gamma_t = \frac{\nu}{C} g_t$$

Heuristic:

$$\gamma_t = \frac{g_t}{C + 2 \frac{\|M\|^2 D^2}{\mu}}$$

$$\mathcal{L}(x, y) := \frac{\mu}{2} \|x - x^*\|_2^2 + (x - x^*)^T M (y - y^*) - \frac{\mu}{2} \|y - y^*\|_2^2$$

- $\mathcal{X} = \mathcal{Y} := [0, 1]^d$
- $d = 30$
- $C := 2LD^2$
- $L = \mu$

References

- [1] J. H. Hammond. *Solving asymmetric variational inequality problems and systems of equations with generalized nonlinear programming algorithms*. PhD thesis, Massachusetts Institute of Technology, 1984.
- [2] N. He and Z. Harchaoui. Semi-proximal mirror-prox for nonsmooth composite minimization. In *NIPS*, 2015.
- [3] S. Karlin. *Mathematical methods and theory in games, programming and economics*, 1960.
- [4] S. Lacoste-Julien and M. Jaggi. On the global linear convergence of Frank-Wolfe optimization variants. In *NIPS*, 2015.
- [5] J. Robinson. An iterative method of solving a game. *Annals of mathematics*, 1951.