

Frank-Wolfe Algorithms for Saddle Point problems



Gauthier Gidel^{1,3}



Tony Jebara²



Simon Lacoste-Julien³

¹INRIA Paris, Sierra Team

²Department of CS, Columbia University

³Department of CS & OR (DIRO) Université de Montréal

12th July 2017

Overview

- ▶ Frank-Wolfe algorithm (FW) gained in popularity in the last couple of years.
- ▶ Main advantage: FW only needs LMO.
- ▶ Extend FW properties to solve saddle point problems¹.
- ▶ **Straightforward** extension but **Non trivial** analysis.

¹Gauthier Gidel, Tony Jebara, and Simon Lacoste-Julien. “Frank-Wolfe Algorithms for Saddle Point Problems”. In: *AISTATS*. 2017.

Overview

- ▶ Frank-Wolfe algorithm (FW) gained in popularity in the last couple of years.
- ▶ Main advantage: FW only needs LMO.
- ▶ Extend FW properties to solve saddle point problems¹.
- ▶ **Straightforward** extension but **Non trivial** analysis.

Question for the audience: Call for application

¹Gauthier Gidel, Tony Jebara, and Simon Lacoste-Julien. “Frank-Wolfe Algorithms for Saddle Point Problems”. In: *AISTATS*. 2017.

Saddle point and link with variational inequalities

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 (x^*, y^*) is called a *Saddle Point*.

Saddle point and link with variational inequalities

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 (x^*, y^*) is called a *Saddle Point*.

► **Necessary stationary conditions:**

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

Saddle point and link with variational inequalities

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 (x^*, y^*) is called a *Saddle Point*.

► *Necessary stationary conditions:*

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

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

Saddle point and link with variational inequalities

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 (x^*, y^*) is called a *Saddle Point*.

► *Necessary stationary conditions:*

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

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

► *Variational inequality:*

$$\forall z \in \mathcal{X} \times \mathcal{Y} \quad \langle z - z^*, g(z^*) \rangle \geq 0$$

where $(x^*, y^*) = z^*$ and $g(z) = (\nabla_x \mathcal{L}(z), -\nabla_y \mathcal{L}(z))$

Saddle point and link with variational inequalities

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 (x^*, y^*) is called a *Saddle Point*.

- ▶ *Necessary stationary conditions:*

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

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

- ▶ *Variational inequality:*

$$\forall z \in \mathcal{X} \times \mathcal{Y} \quad \langle z - z^*, g(z^*) \rangle \geq 0$$

where $(x^*, y^*) = z^*$ and $g(z) = (\nabla_x \mathcal{L}(z), -\nabla_y \mathcal{L}(z))$

- ▶ *Sufficient condition: 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.

Motivations: games and robust learning

- ▶ *Zero-sum games with two players:*

$$\min_{\mathbf{x} \in \Delta(I)} \max_{\mathbf{y} \in \Delta(J)} \mathbf{x}^\top M \mathbf{y}$$

²J. Wen, C. Yu, and R. Greiner. “Robust Learning under Uncertain Test Distributions: Relating Covariate Shift to Model Misspecification.” In: *ICML*. 2014.

Motivations: games and robust learning

- ▶ *Zero-sum games with two players:*

$$\min_{x \in \Delta(I)} \max_{y \in \Delta(J)} x^\top M y$$

- ▶ *Robust learning:*² We want to learn

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(f_\theta(x_i), y_i) + \lambda \Omega(\theta)$$

with an uncertainty regarding the data:

$$\min_{\theta \in \Theta} \max_{w \in \Delta_n} \sum_{i=1}^n \omega_i \ell(f_\theta(x_i), y_i) + \lambda \Omega(\theta)$$

Minimize the **worst case** → gives **robustness**

²J. Wen, C. Yu, and R. Greiner. “Robust Learning under Uncertain Test Distributions: Relating Covariate Shift to Model Misspecification.” In: *ICML*. 2014.

Problem with Hard projection

The *structured SVM*:

$$\min_{\omega \in \mathbb{R}^d} \lambda \Omega(\omega) + \underbrace{\frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}_i} (L_i(y) - \langle \omega, \phi_i(y) \rangle)}_{\text{structured empirical loss}}$$

Problem with Hard projection

The *structured SVM*:

$$\min_{\omega \in \mathbb{R}^d} \lambda \Omega(\omega) + \underbrace{\frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}_i} (L_i(y) - \langle \omega, \phi_i(y) \rangle)}_{\text{structured empirical loss}}$$

Regularization: penalized \rightarrow constrained.

$$\min_{\Omega(\omega) \leq \beta} \max_{\alpha \in \Delta(|\mathcal{Y}|)} b^T \alpha - \omega^T M \alpha$$

Problem with Hard projection

The *structured SVM*:

$$\min_{\omega \in \mathbb{R}^d} \lambda \Omega(\omega) + \underbrace{\frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}_i} (L_i(y) - \langle \omega, \phi_i(y) \rangle)}_{\text{structured empirical loss}}$$

Regularization: penalized \rightarrow constrained.

$$\min_{\Omega(\omega) \leq \beta} \max_{\alpha \in \Delta(|\mathcal{Y}|)} b^T \alpha - \omega^T M \alpha$$

Problem with Hard projection

The *structured SVM*:

$$\min_{\omega \in \mathbb{R}^d} \lambda \Omega(\omega) + \underbrace{\frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}_i} (L_i(y) - \langle \omega, \phi_i(y) \rangle)}_{\text{structured empirical loss}}$$

Regularization: penalized \rightarrow constrained.

$$\min_{\Omega(\omega) \leq \beta} \max_{\alpha \in \Delta(|\mathcal{Y}|)} b^T \alpha - \omega^T M \alpha$$

Problem with Hard projection

The *structured SVM*:

$$\min_{\omega \in \mathbb{R}^d} \lambda \Omega(\omega) + \underbrace{\frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}_i} (L_i(y) - \langle \omega, \phi_i(y) \rangle)}_{\text{structured empirical loss}}$$

Regularization: penalized \rightarrow constrained.

$$\min_{\Omega(\omega) \leq \beta} \max_{\alpha \in \Delta(|\mathcal{Y}|)} b^T \alpha - \omega^T M \alpha$$

Difficult to project when:

Problem with Hard projection

The *structured SVM*:

$$\min_{\omega \in \mathbb{R}^d} \lambda \Omega(\omega) + \underbrace{\frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}_i} (L_i(y) - \langle \omega, \phi_i(y) \rangle)}_{\text{structured empirical loss}}$$

Regularization: penalized \rightarrow constrained.

$$\min_{\Omega(\omega) \leq \beta} \max_{\alpha \in \Delta(|\mathcal{Y}|)} b^T \alpha - \omega^T M \alpha$$

Difficult to project when:

- ▶ **Structured sparsity** norm (group lasso norm).

Problem with Hard projection

The *structured SVM*:

$$\min_{\omega \in \mathbb{R}^d} \lambda \Omega(\omega) + \underbrace{\frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}_i} (L_i(y) - \langle \omega, \phi_i(y) \rangle)}_{\text{structured empirical loss}}$$

Regularization: penalized \rightarrow constrained.

$$\min_{\Omega(\omega) \leq \beta} \max_{\alpha \in \Delta(|\mathcal{Y}|)} b^T \alpha - \omega^T M \alpha$$

Difficult to project when:

- ▶ **Structured sparsity** norm (group lasso norm).
- ▶ The output \mathcal{Y} is structured: **exponential size**.

Standard approaches in literature

- ▶ Projected gradient algorithm.

$$\mathbf{x}^{(t+1)} = P_{\mathcal{X}}(\mathbf{x}^{(t)} - \eta \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}))$$

$$\mathbf{y}^{(t+1)} = P_{\mathcal{Y}}(\mathbf{y}^{(t)} + \eta \nabla_{\mathbf{y}} \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}))$$

³GM Korpelevich. “The extragradient method for finding saddle points and other problems”. In: *Matecon* (1976).

Standard approaches in literature

- ▶ Projected gradient algorithm.

$$\mathbf{x}^{(t+1)} = P_{\mathcal{X}}(\mathbf{x}^{(t)} - \eta \nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}))$$

$$\mathbf{y}^{(t+1)} = P_{\mathcal{Y}}(\mathbf{y}^{(t)} + \eta \nabla_y \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}))$$

- ▶ Projected extra-gradient³.

$$\bar{\mathbf{x}}^{(t+1)} = P_{\mathcal{X}}(\mathbf{x}^{(t)} - \eta \nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}))$$

$$\bar{\mathbf{y}}^{(t+1)} = P_{\mathcal{Y}}(\mathbf{y}^{(t)} + \eta \nabla_y \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}))$$

Intuition: *lookahead* move: look at what your opponent *would* do before deciding your move.

$$\mathbf{x}^{(t+1)} = P_{\mathcal{X}}(\mathbf{x}^{(t)} - \eta \nabla_x \mathcal{L}(\bar{\mathbf{x}}^{(t+1)}, \bar{\mathbf{y}}^{(t+1)}))$$

$$\mathbf{y}^{(t+1)} = P_{\mathcal{Y}}(\mathbf{y}^{(t)} + \eta \nabla_y \mathcal{L}(\bar{\mathbf{x}}^{(t+1)}, \bar{\mathbf{y}}^{(t+1)}))$$

Prevents oscillations for non strongly convex objective.

³GM Korpelevich. “The extragradient method for finding saddle points and other problems”. In: *Matecon* (1976).

Standard approaches in literature

- ▶ Gradient method works for non-smooth optimization, but

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

⁴N. He and Z. Harchaoui. “Semi-proximal Mirror-Prox for Nonsmooth Composite Minimization”. In: *NIPS*. 2015.

Standard approaches in literature

- ▶ Gradient method works for non-smooth optimization, but

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

- ▶ Extragradient method works for smooth optimization,

$$(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \rightarrow (\mathbf{x}^*, \mathbf{y}^*)$$

⁴N. He and Z. Harchaoui. “Semi-proximal Mirror-Prox for Nonsmooth Composite Minimization”. In: *NIPS*. 2015.

Standard approaches in literature

- ▶ Gradient method works for non-smooth optimization, but

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

- ▶ Extragradient method works for smooth optimization,

$$(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \rightarrow (\mathbf{x}^*, \mathbf{y}^*)$$

Even when projections are expensive:

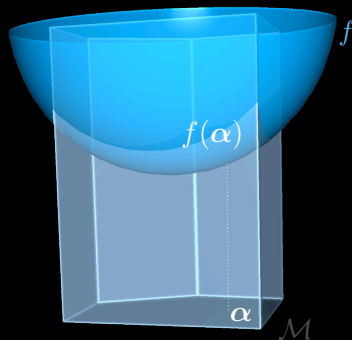
Can use LMO to compute approximate projections⁴.

⁴N. He and Z. Harchaoui. “Semi-proximal Mirror-Prox for Nonsmooth Composite Minimization”. In: *NIPS*. 2015.

The FW algorithm

Algorithm Frank-Wolfe algorithm

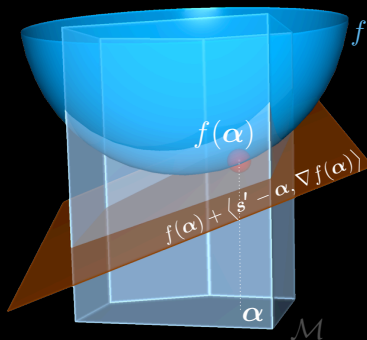
- 1: Let $\mathbf{x}^{(0)} \in \mathcal{X}$
 - 2: **for** $t = 0 \dots T$ **do**
 - 3: Compute $\mathbf{r}^{(t)} = \nabla f(\mathbf{x}^{(t)})$
 - 4: Compute $\mathbf{s}^{(t)} \in \underset{\mathbf{s} \in \mathcal{X}}{\operatorname{argmin}} \langle \mathbf{s}, \mathbf{r}^{(t)} \rangle$
 - 5: Compute $g_t := \langle \mathbf{x}^{(t)} - \mathbf{s}^{(t)}, \mathbf{r}^{(t)} \rangle$
 - 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**
-



The FW algorithm

Algorithm Frank-Wolfe algorithm

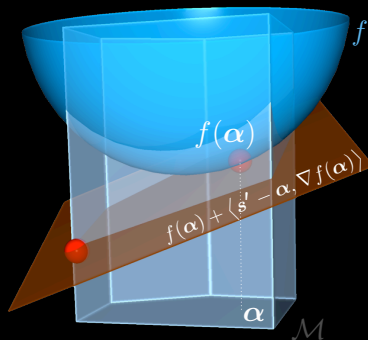
- 1: Let $\mathbf{x}^{(0)} \in \mathcal{X}$
 - 2: **for** $t = 0 \dots T$ **do**
 - 3: **Compute** $\mathbf{r}^{(t)} = \nabla f(\mathbf{x}^{(t)})$
 - 4: **Compute** $\mathbf{s}^{(t)} \in \operatorname{argmin}_{\mathbf{s} \in \mathcal{X}} \langle \mathbf{s}, \mathbf{r}^{(t)} \rangle$
 - 5: **Compute** $g_t := \langle \mathbf{x}^{(t)} - \mathbf{s}^{(t)}, \mathbf{r}^{(t)} \rangle$
 - 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**
-



The FW algorithm

Algorithm Frank-Wolfe algorithm

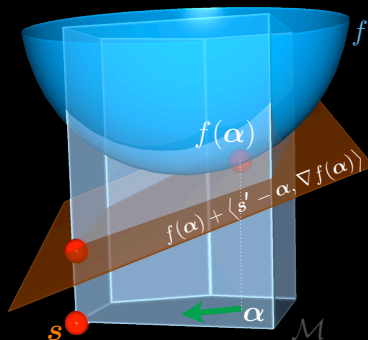
- 1: Let $\mathbf{x}^{(0)} \in \mathcal{X}$
 - 2: **for** $t = 0 \dots T$ **do**
 - 3: Compute $\mathbf{r}^{(t)} = \nabla f(\mathbf{x}^{(t)})$
 - 4: **Compute** $\mathbf{s}^{(t)} \in \underset{\mathbf{s} \in \mathcal{X}}{\operatorname{argmin}} \langle \mathbf{s}, \mathbf{r}^{(t)} \rangle$
 - 5: Compute $g_t := \langle \mathbf{x}^{(t)} - \mathbf{s}^{(t)}, \mathbf{r}^{(t)} \rangle$
 - 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**
-



The FW algorithm

Algorithm 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 $\mathbf{s}^{(t)} \in \operatorname{argmin}_{\mathbf{s} \in \mathcal{X}} \langle \mathbf{s}, \mathbf{r}^{(t)} \rangle$
 - 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 Saddle point FW algorithm

1: **for** $t = 0 \dots T$ **do**

2: **Compute** $\mathbf{r}^{(t)} := \begin{pmatrix} \nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \\ -\nabla_y \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \end{pmatrix}$

3: **Compute** $\mathbf{s}^{(t)} \in \operatorname{argmin}_{\mathbf{z} \in \mathcal{X} \times \mathcal{Y}} \langle \mathbf{z}, \mathbf{r}^{(t)} \rangle$

4: **Compute** $g_t := \langle \mathbf{z}^{(t)} - \mathbf{s}^{(t)}, \mathbf{r}^{(t)} \rangle$

5: **if** $g_t \leq \epsilon$ **then return** $\mathbf{z}^{(t)}$

6: **Let** $\gamma = \min(1, \frac{\nu}{C} g_t)$ **or** $\gamma = \frac{2}{2+t}$

7: **Update** $\mathbf{z}^{(t+1)} := (1 - \gamma)\mathbf{z}^{(t)} + \gamma\mathbf{s}^{(t)}$

8: **end for**

- ▶ Originally proposed by Hammond⁴ with $\gamma_t = 1/(t+1)$.

⁵J. Hammond. “Solving asymmetric variational inequality problems and systems of equations with generalized nonlinear programming algorithms”. PhD thesis. MIT, 1984.

Algorithm Saddle point FW algorithm

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

- ▶ Originally proposed by Hammond⁴ with $\gamma_t = 1/(t+1)$.

⁵J. Hammond. “Solving asymmetric variational inequality problems and systems of equations with generalized nonlinear programming algorithms”. PhD thesis. MIT, 1984.

Algorithm Saddle point FW algorithm

```

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

```

- Originally proposed by Hammond⁴ with $\gamma_t = 1/(t+1)$.

⁵J. Hammond. “Solving asymmetric variational inequality problems and systems of equations with generalized nonlinear programming algorithms”. PhD thesis. MIT, 1984.

Algorithm Saddle point FW algorithm

1: **for** $t = 0 \dots T$ **do**

2: Compute $\mathbf{r}^{(t)} := \begin{pmatrix} \nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \\ -\nabla_y \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \end{pmatrix}$

3: Compute $\mathbf{s}^{(t)} \in \operatorname{argmin}_{\mathbf{z} \in \mathcal{X} \times \mathcal{Y}} \langle \mathbf{z}, \mathbf{r}^{(t)} \rangle$

4: Compute $g_t := \langle \mathbf{z}^{(t)} - \mathbf{s}^{(t)}, \mathbf{r}^{(t)} \rangle$

5: **if** $g_t \leq \epsilon$ **then return** $\mathbf{z}^{(t)}$

6: Let $\gamma = \min(1, \frac{\nu}{C} g_t)$ **or** $\gamma = \frac{2}{2+t}$

7: Update $\mathbf{z}^{(t+1)} := (1 - \gamma)\mathbf{z}^{(t)} + \gamma\mathbf{s}^{(t)}$

8: **end for**

- ▶ Originally proposed by Hammond⁴ with $\gamma_t = 1/(t+1)$.

⁵J. Hammond. “Solving asymmetric variational inequality problems and systems of equations with generalized nonlinear programming algorithms”. PhD thesis. MIT, 1984.

Algorithm Saddle point FW algorithm

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

- ▶ Originally proposed by Hammond⁴ with $\gamma_t = 1/(t+1)$.

⁵J. Hammond. “Solving asymmetric variational inequality problems and systems of equations with generalized nonlinear programming algorithms”. PhD thesis. MIT, 1984.

Algorithm Saddle point FW algorithm

1: **for** $t = 0 \dots T$ **do**

2: Compute $\mathbf{r}^{(t)} := \begin{pmatrix} \nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \\ -\nabla_y \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \end{pmatrix}$

3: Compute $\mathbf{s}^{(t)} \in \underset{\mathbf{z} \in \mathcal{X} \times \mathcal{Y}}{\operatorname{argmin}} \langle \mathbf{z}, \mathbf{r}^{(t)} \rangle$

4: Compute $g_t := \langle \mathbf{z}^{(t)} - \mathbf{s}^{(t)}, \mathbf{r}^{(t)} \rangle$

5: **if** $g_t \leq \epsilon$ **then return** $\mathbf{z}^{(t)}$

6: Let $\gamma = \min(1, \frac{\nu}{C} g_t)$ **or** $\gamma = \frac{2}{2+t}$

7: Update $\mathbf{z}^{(t+1)} := (1 - \gamma)\mathbf{z}^{(t)} + \gamma\mathbf{s}^{(t)}$

8: **end for**

- ▶ Originally proposed by Hammond⁴ with $\gamma_t = 1/(t+1)$.
- ▶ One can define FW extension with **away** step.

⁵J. Hammond. “Solving asymmetric variational inequality problems and systems of equations with generalized nonlinear programming algorithms”. PhD thesis. MIT, 1984.

Algorithm Saddle point FW algorithm

1: **for** $t = 0 \dots T$ **do**

2: Compute $\mathbf{r}^{(t)} := \begin{pmatrix} \nabla_x \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \\ -\nabla_y \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \end{pmatrix}$

3: Compute $\mathbf{s}^{(t)} \in \underset{\mathbf{z} \in \mathcal{X} \times \mathcal{Y}}{\operatorname{argmin}} \langle \mathbf{z}, \mathbf{r}^{(t)} \rangle$

4: Compute $g_t := \langle \mathbf{z}^{(t)} - \mathbf{s}^{(t)}, \mathbf{r}^{(t)} \rangle$

5: **if** $g_t \leq \epsilon$ **then return** $\mathbf{z}^{(t)}$

6: Let $\gamma = \min(1, \frac{\nu}{C} g_t)$ **or** $\gamma = \frac{2}{2+t}$

7: Update $\mathbf{z}^{(t+1)} := (1 - \gamma)\mathbf{z}^{(t)} + \gamma\mathbf{s}^{(t)}$

8: **end for**

- ▶ Originally proposed by Hammond⁴ with $\gamma_t = 1/(t+1)$.
- ▶ One can define FW extension with **away** step.
- ▶ Crucial for our linear convergence results.

⁵J. Hammond. “Solving asymmetric variational inequality problems and systems of equations with generalized nonlinear programming algorithms”. PhD thesis. MIT, 1984.

Algorithm Saddle point FW algorithm

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

► $\gamma_t = \frac{1}{1+t} \Rightarrow \mathbf{z}^{(t)} = \frac{1}{t} \sum_{i=0}^{t-1} \mathbf{s}^{(i)}$.

- Originally proposed by Hammond⁴ with $\gamma_t = 1/(t+1)$.
- One can define FW extension with **away** step.
- Crucial for our linear convergence results.

⁵J. Hammond. “Solving asymmetric variational inequality problems and systems of equations with generalized nonlinear programming algorithms”. PhD thesis. MIT, 1984.

Algorithm Saddle point FW algorithm

```

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

```

▶ $\gamma_t = \frac{1}{1+t} \Rightarrow \mathbf{z}^{(t)} = \frac{1}{t} \sum_{i=0}^{t-1} \mathbf{s}^{(i)}$.

▶ $(\gamma_t = \frac{1}{1+t}) + \text{Bilinear objective} \leftrightarrow \text{fictitious play algorithm.}^5$

- ▶ Originally proposed by Hammond⁴ with $\gamma_t = 1/(t+1)$.
- ▶ One can define FW extension with **away** step.
- ▶ Crucial for our linear convergence results.

⁵J. Hammond. “Solving asymmetric variational inequality problems and systems of equations with generalized nonlinear programming algorithms”. PhD thesis. MIT, 1984.

Advantages of SP-FW

Same main property as FW:

Only LMO (linear minimization oracle).

Advantages of SP-FW

Same main property as FW:

Only LMO (linear minimization oracle).

Same other **advantages** as FW:

- ▶ Convergence certificate g_t for free.

Advantages of SP-FW

Same main property as FW:

Only LMO (linear minimization oracle).

Same other **advantages** as FW:

- ▶ Convergence certificate g_t for free.
- ▶ Affine invariance of the algorithm.

Advantages of SP-FW

Same main property as FW:

Only LMO (linear minimization oracle).

Same other **advantages** as FW:

- ▶ Convergence certificate g_t for free.
- ▶ Affine invariance of the algorithm.
- ▶ *Sparsity* of the iterates.

Advantages of SP-FW

Same main property as FW:

Only LMO (linear minimization oracle).

Same other **advantages** as FW:

- ▶ Convergence certificate g_t for free.
- ▶ Affine invariance of the algorithm.
- ▶ *Sparsity* of the iterates.
- ▶ Universal step size $\gamma_t := \frac{2}{2+t}$, adaptive step size $\gamma_t := \frac{\nu}{C} g_t$.

Advantages of SP-FW

Same main property as FW:

Only LMO (linear minimization oracle).

Same other **advantages** as FW:

- ▶ Convergence certificate g_t for free.
- ▶ Affine invariance of the algorithm.
- ▶ *Sparsity* of the iterates.
- ▶ Universal step size $\gamma_t := \frac{2}{2+t}$, adaptive step size $\gamma_t := \frac{\nu}{C} g_t$.

Main **difference** with FW:

- ▶ **No** line-search.

Advantages of SP-FW

Same main property as FW:

Only LMO (linear minimization oracle).

Same other **advantages** as FW:

- ▶ Convergence certificate g_t for free.
- ▶ Affine invariance of the algorithm.
- ▶ *Sparsity* of the iterates.
- ▶ Universal step size $\gamma_t := \frac{2}{2+t}$, adaptive step size $\gamma_t := \frac{\nu}{C} g_t$.

Main **difference** with FW:

- ▶ **No** line-search.

When constraint set is a “complicated” *structured* polytope:
projection is *difficult* whereas LMO is *tractable*.

Hypothesis

Similar hypothesis as AFW:

- ▶ \mathcal{L} is L -smooth and μ -strongly convex-concave.
- ▶ \mathcal{X} and \mathcal{Y} polytopes.

Hypothesis

Similar hypothesis as AFW:

- ▶ \mathcal{L} is L -smooth and μ -strongly convex-concave.
- ▶ \mathcal{X} and \mathcal{Y} polytopes.
- ▶ Additional assumption on **bilinearity**:

$$\mathcal{L}(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{x}^\top M \mathbf{y} - h(\mathbf{y})$$

Roughly, $\|M\|$ smaller than the strong convexity constant.

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

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

Theoretical contribution

SP extension of FW with *away step*⁶:

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

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

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

$$\min_{s \leq t} g_s \leq O\left(\frac{1}{t}\right)$$

- $k(t)$: number of non drop steps, $k(t) \geq t/3$.

⁶Gauthier Gidel, Tony Jebara, and Simon Lacoste-Julien. “Frank-Wolfe Algorithms for Saddle Point Problems”. In: *AISTATS*. 2017.

Theoretical contribution

SP extension of FW with *away step*⁶:

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

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

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

$$\min_{s \leq t} g_s \leq O\left(\frac{1}{t}\right)$$

- ▶ $k(t)$: number of non drop steps, $k(t) \geq t/3$.
- ▶ Proof use recent advances on AFW \rightarrow **growth condition**.

⁶Gauthier Gidel, Tony Jebara, and Simon Lacoste-Julien. “Frank-Wolfe Algorithms for Saddle Point Problems”. In: *AISTATS*. 2017.

Theoretical contribution

SP extension of FW with *away step*⁷:

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

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

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

$$\min_{s \leq t} g_s \leq O\left(\frac{1}{t}\right)$$

- ▶ $k(t)$: number of non drop steps, $k(t) \geq t/3$.
- ▶ Proof use recent advances on AFW \rightarrow **growth condition**.
- ▶ Partially answering a **30 years old conjecture**⁸.
 - strongly monotone obj with step size $\frac{1}{t+1}$ over polytope.

⁷Gauthier Gidel, Tony Jebara, and Simon Lacoste-Julien. “Frank-Wolfe Algorithms for Saddle Point Problems”. In: *AISTATS*. 2017.

⁸J. Hammond. “Solving asymmetric variational inequality problems and systems of equations with generalized nonlinear programming algorithms”. PhD thesis. MIT, 1984.

Difficulties for saddle point

Usual **descent Lemma**:

$$h_{t+1} \leq h_t - \underbrace{\gamma_t g_t}_{\geq 0} + \gamma_t^2 \frac{L \|d^{(t)}\|^2}{2}$$

With γ_t small enough the sequence **decreases**.

Difficulties for saddle point

Usual **descent Lemma**:

$$h_{t+1} \leq h_t - \underbrace{\gamma_t g_t}_{\geq 0} + \gamma_t^2 \frac{L \|d^{(t)}\|^2}{2}$$

With γ_t small enough the sequence **decreases**.

For saddle point problem the Lipschitz gradient property gives

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

- ▶ Cannot control the **oscillation** of the sequence.
- ▶ Must introduce other quantities to establish convergence.

Difficulties for saddle point

Standard merit functions: *primal + dual gaps*

$$h_t := \max_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} \mathcal{L}(\mathbf{x}, \mathbf{y}^{(t)}) \geq 0.$$

⁹Gauthier Gidel, Tony Jebara, and Simon Lacoste-Julien. “Frank-Wolfe Algorithms for Saddle Point Problems”. In: *AISTATS*. 2017.

Difficulties for saddle point

Standard merit functions: *primal + dual gaps*

$$h_t := \max_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} \mathcal{L}(\mathbf{x}, \mathbf{y}^{(t)}) \geq 0.$$

Problem: $\hat{\mathbf{y}}^{(t)} := \arg \max_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y})$ depends on t .

⁹Gauthier Gidel, Tony Jebara, and Simon Lacoste-Julien. “Frank-Wolfe Algorithms for Saddle Point Problems”. In: *AISTATS*. 2017.

Difficulties for saddle point

Standard merit functions: *primal + dual gaps*

$$h_t := \max_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} \mathcal{L}(\mathbf{x}, \mathbf{y}^{(t)}) \geq 0.$$

Problem: $\hat{\mathbf{y}}^{(t)} := \arg \max_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y})$ depends on t .

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

We have,

$$0 \leq w_t \leq h_t \leq g_t$$

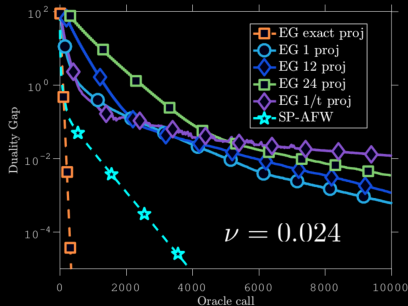
In general, w_t can be zero even if we have not reached a solution. But for strongly convex-concave function⁹

$$h_t \leq Cte\sqrt{w_t}$$

⁹Gauthier Gidel, Tony Jebara, and Simon Lacoste-Julien. “Frank-Wolfe Algorithms for Saddle Point Problems”. In: *AISTATS*. 2017.

Toy Experiments

- ▶ SP-AFW vs. Extragradient with approximate projection.



Theoretical step-size

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

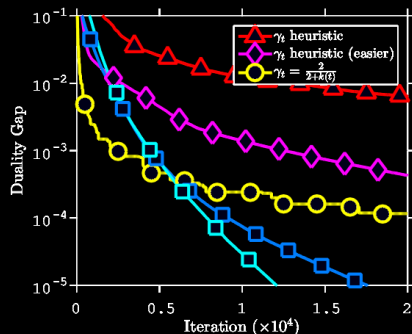
EG : [He & Harchaoui NIPS 2015]

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

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

Toy Experiments

- ▶ SP-AFW with heuristic step-size. (When $\nu < 0$)



Heuristic step-size.

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

Recall: theoretical step-size

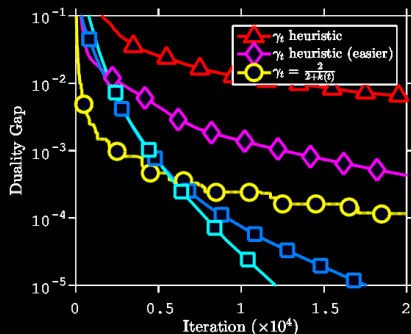
$$\gamma_t = \frac{\nu}{C} gt.$$

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

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

Toy Experiments

- ▶ SP-AFW with heuristic step-size. (When $\nu < 0$)



Heuristic step-size.

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

Recall: theoretical step-size

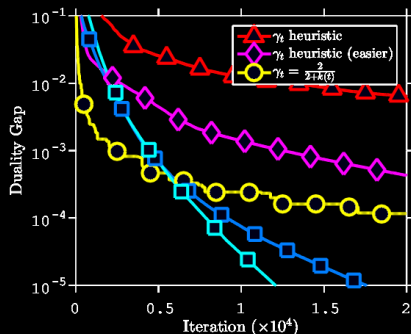
$$\gamma_t = \frac{\nu}{C} gt.$$

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

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

Toy Experiments

- ▶ SP-AFW with heuristic step-size. (When $\nu < 0$)



Heuristic step-size.

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

Recall: theoretical step-size

$$\gamma_t = \frac{\nu}{C} gt.$$

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

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

Conclusion

- ▶ SP-FW one of the first SP solver only working with LMO.

Conclusion

- ▶ SP-FW one of the first SP solver only working with LMO.
- ▶ FW resurgence lead to new *structured* problems.

Conclusion

- ▶ SP-FW one of the first SP solver only working with LMO.
- ▶ FW resurgence lead to new *structured* problems.
- ▶ Same hope as FW for SP-FW

Conclusion

- ▶ SP-FW one of the first SP solver only working with LMO.
- ▶ FW resurgence lead to new *structured* problems.
- ▶ Same hope as FW for SP-FW ∇

Call for applications !

Conclusion

- ▶ SP-FW one of the first SP solver only working with LMO.
- ▶ FW resurgence lead to new *structured* problems.
- ▶ Same hope as FW for SP-FW ∇

Call for applications !

- ▶ With a bilinear objective this algorithm is *highly related* to the *fictitious play algorithm*.

Conclusion

- ▶ SP-FW one of the first SP solver only working with LMO.
- ▶ FW resurgence lead to new *structured* problems.
- ▶ Same hope as FW for SP-FW ∇

Call for applications !

- ▶ With a bilinear objective this algorithm is *highly related* to the *fictitious play algorithm*.
- ▶ Rich interplay tapping into this game theory literature.

Conclusion

- ▶ SP-FW one of the first SP solver only working with LMO.
- ▶ FW resurgence lead to new *structured* problems.
- ▶ Same hope as FW for SP-FW ∇

Call for applications !

- ▶ With a bilinear objective this algorithm is *highly related* to the *fictitious play algorithm*.
- ▶ Rich interplay tapping into this game theory literature.
- ▶ Still many theoretical opened questions.
 - ↳ Karlin's conjecture.¹⁰
 - ↳ Convergence without assumption on the bilinearity.

¹⁰Samuel Karlin. *Mathematical methods and theory in games, programming and economics*. 1960.

Thank You !

Slides available on www.di.ens.fr/~gidel.

Problems with difficult projection

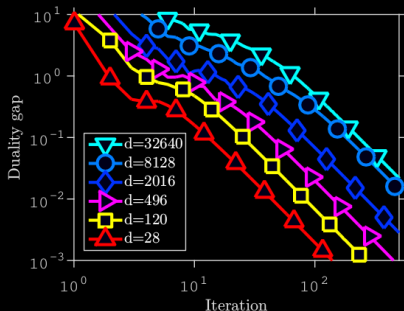
University game:

1. Game between two universities (A and B).
2. Admitting d students and have to assign pairs of students into dorms.
3. The game has a payoff matrix M belonging to $\mathbb{R}^{(d(d-1)/2)^2}$.
4. $M_{ij,kl}$ 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 .
5. Here the actions are both in the *marginal polytope* of all perfect *unipartite matchings*.

Hard to project on this polytope whereas the LMO can be solved efficiently with the blossom algorithm¹¹.

¹¹J. Edmonds. “Paths, trees and flowers”. In: *Canadian Journal of Mathematics* (1965).

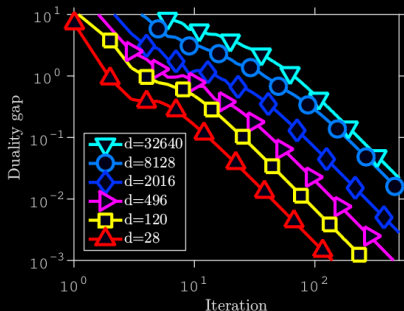
Experiments



- Sublinear convergence rate (faster than expected $O(t^{-2})$)

Figure: SP-FW on the University game.

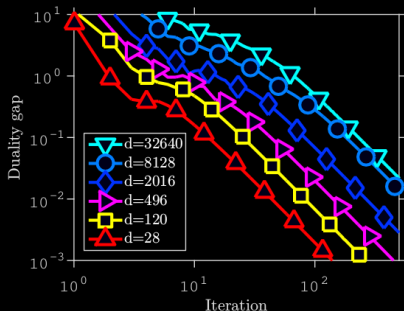
Experiments



- ▶ Sublinear convergence rate (faster than expected $O(t^{-2})$)
- ▶ Best theoretical rate proved: $O(t^{-1/d})$

Figure: SP-FW on the University game.

Experiments



- ▶ Sublinear convergence rate (faster than expected $O(t^{-2})$)
- ▶ Best theoretical rate proved: $O(t^{-1/d})$
- ▶ Scale well with dimension.

Figure: SP-FW on the University game.