# Frank-Wolfe Methods for Optimization and Machine Learning

Cyrille W. Combettes

School of Industrial and Systems Engineering Georgia Institute of Technology

April 16, 2021



#### 1 Introduction

- 2 The Frank-Wolfe algorithm
- **3** Boosting Frank-Wolfe by chasing gradients
- **4** Adaptive Frank-Wolfe for large-scale optimization
- 5 The approximate Carathéodory problem

#### Introduction

We consider

min 
$$f(x)$$
  
s.t.  $x \in C$ 

where

- $\mathcal{C} \subset \mathbb{R}^n$  is a compact convex set
- $f: \mathbb{R}^n \to \mathbb{R}$  is a smooth convex function

We consider

min 
$$f(x)$$
  
s.t.  $x \in C$ 

where

- $\mathcal{C} \subset \mathbb{R}^n$  is a compact convex set
- $f: \mathbb{R}^n \to \mathbb{R}$  is a smooth convex function

#### Example

Sparse logistic regression

• Low-rank matrix completion

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{m} \sum_{i=1}^m \ln(1 + \exp(-y_i \langle \mathbf{a}_i, \mathbf{x} \rangle))$$
s.t.  $\|\mathbf{x}\|_1 \leqslant \tau$ 

 $\min_{X \in \mathbb{R}^{m \times n}} \frac{1}{2|\mathcal{I}|} \sum_{(i,j) \in \mathcal{I}} (Y_{i,j} - X_{i,j})^2$ s.t.  $\|X\|_{\text{nuc}} \leq \tau$  • A natural approach is to use any efficient method and add projections back onto  ${\mathcal C}$  to ensure feasibility

- A natural approach is to use any efficient method and add projections back onto  ${\cal C}$  to ensure feasibility



- A natural approach is to use any efficient method and add projections back onto  ${\cal C}$  to ensure feasibility
- However, in many situations projections onto  $\mathcal C$  are very expensive

- A natural approach is to use any efficient method and add projections back onto  ${\cal C}$  to ensure feasibility
- However, in many situations projections onto  ${\mathcal C}$  are very expensive
- This is an issue with the method of projections, not necessarily with the geometry of C: linear minimizations over C can still be relatively cheap

- A natural approach is to use any efficient method and add projections back onto  ${\cal C}$  to ensure feasibility
- However, in many situations projections onto  $\ensuremath{\mathcal{C}}$  are very expensive
- This is an issue with the method of projections, not necessarily with the geometry of C: linear minimizations over C can still be relatively cheap
- We compare

$$rgmin_{x\in\mathcal{C}} \langle x,y 
angle$$
 and  $rgmin_{x\in\mathcal{C}} \|x-y\|$ 

on several sets commonly used in optimization

Set $C$	Linear minimization	Projection
$\ell_1/\ell_2/\ell_{\infty}\text{-ball}$	$\mathcal{O}(n)$ $\mathcal{O}(n)$	$\mathcal{O}(n)$ $\mathcal{O}(no^{2}  _{V} - x^{*}  _{2}^{2}/c^{2})$
Nuclear norm-ball	$\mathcal{O}(\nu \ln(m+n)\sqrt{\sigma_1}/\sqrt{\varepsilon})$	$\mathcal{O}(mp   y - x   _2/\varepsilon)$ $\mathcal{O}(mn \min\{m, n\})$
Flow polytope Birkhoff polytope	$\mathcal{O}(m+n)$ $\mathcal{O}(n^3)$	$\mathcal{O}(m^3n + n^2)$ $\mathcal{O}(n^2d_r^2/\varepsilon^2)$
Permutahedron	$\mathcal{O}(n\ln(n))$	$\mathcal{O}(n\ln(n) + n)$

Set $C$	Linear minimization	Projection
$\begin{array}{l} \ell_1/\ell_2/\ell_{\infty}\text{-ball} \\ \ell_p\text{-ball}, \ p \in ]1, \infty[\setminus\{2\} \\ \text{Nuclear norm-ball} \\ \text{Flow polytope} \\ \text{Birkhoff polytope} \end{array}$	$ \begin{array}{c} \mathcal{O}(n) \\ \mathcal{O}(n) \\ \mathcal{O}(\nu \ln(m+n)\sqrt{\sigma_1}/\sqrt{\varepsilon}) \\ \mathcal{O}(m+n) \\ \mathcal{O}(n^3) \end{array} $	$\mathcal{O}(n) \\ \mathcal{O}(n\rho^2    y - x^*   _2^2 / \varepsilon^2) \\ \mathcal{O}(mn\min\{m, n\}) \\ \tilde{\mathcal{O}}(m^3 n + n^2) \\ \mathcal{O}(n^2 d_z^2 / \varepsilon^2)$
Permutahedron	$\mathcal{O}(n\ln(n))$	$\mathcal{O}(n\ln(n)+n)$

Set $\mathcal{C}$	Linear minimization	Projection
$\begin{array}{l} \ell_1/\ell_2/\ell_\infty\text{-ball}\\ \ell_p\text{-ball}, \ p\in ]1, \infty[\backslash\{2\}\\ \text{Nuclear norm-ball}\\ \text{Flow polytope}\\ \text{Birkhoff polytope} \end{array}$	$ \begin{array}{l} \mathcal{O}(n) \\ \mathcal{O}(n) \\ \mathcal{O}(\nu \ln(m+n)\sqrt{\sigma_1}/\sqrt{\varepsilon}) \\ \mathcal{O}(m+n) \\ \mathcal{O}(n^3) \\ \mathcal{O}(n) \end{array} $	$ \begin{array}{l} \mathcal{O}(n) \\ \mathcal{O}(n\rho^2 \  y - x^* \ _2^2 / \varepsilon^2) \\ \mathcal{O}(mn\min\{m, n\}) \\ \tilde{\mathcal{O}}(m^3 n + n^2) \\ \mathcal{O}(n^2 d_z^2 / \varepsilon^2) \\ \end{array} $
Permutahedron	$O(n \ln(n))$	$O(n \ln(n) + n)$

Example: the  $\ell_1\text{-ball}$ 



Set $C$	Linear minimization	Projection
$\ell_1/\ell_2/\ell_{\infty}\text{-ball}$	$\mathcal{O}(n)$	$\mathcal{O}(n)$ $\mathcal{O}(ne^{2}\ y - x^{*}\ ^{2}/c^{2})$
Nuclear norm-ball	$\mathcal{O}(n)$ $\mathcal{O}(\nu \ln(m+n)\sqrt{\sigma_1}/\sqrt{\varepsilon})$	$\mathcal{O}(mp    y - x   _2/\varepsilon)$ $\mathcal{O}(mn \min\{m, n\})$
Flow polytope Birkhoff polytope	$\mathcal{O}(m+n)$ $\mathcal{O}(n^3)$	$\mathcal{O}(m^3n + n^2)$ $\mathcal{O}(n^2d_r^2/\varepsilon^2)$
Permutahedron	$\mathcal{O}(n\ln(n))$	$\mathcal{O}(n \ln(n) + n)$

	n
$\begin{array}{ccc} \ell_1/\ell_2/\ell_{\infty}\text{-ball} & \mathcal{O}(n) & \mathcal{O}(n) \\ \ell_p\text{-ball}, \ p \in ]1, \infty[\backslash \{2\} & \mathcal{O}(n) & \mathcal{O}(n\rho^2 \  y) \\ \text{Nuclear norm-ball} & \mathcal{O}(\nu \ln(m+n)\sqrt{\sigma_1}/\sqrt{\varepsilon}) & \mathcal{O}(mn \min plus polytope) \\ \text{Flow polytope} & \mathcal{O}(m+n) & \tilde{\mathcal{O}}(m^3 n+1) \\ \text{Birkhoff polytope} & \mathcal{O}(n^3) & \mathcal{O}(n^2 d_z^2/\varepsilon) \\ \text{Permutabedron} & \mathcal{O}(n \ln(n)) & \mathcal{O}(n \ln(n)) \end{array}$	$-x^* \ _2^2 / \varepsilon^2 $ $n\{m, n\} $ $n^2 $ $(x^2)$

Set $C$	Linear minimization	Projection
$\begin{array}{c} \ell_1/\ell_2/\ell_\infty\text{-ball} \\ \ell_p\text{-ball}, \ p \in ]1, \infty[\setminus\{2\} \\ \text{Nuclear norm-ball} \end{array}$	$\mathcal{O}(n) \\ \mathcal{O}(n) \\ \mathcal{O}(\nu \ln(m+n)\sqrt{\sigma_1}/\sqrt{\varepsilon})$	$\mathcal{O}(n) \\ \mathcal{O}(n\rho^2    y - x^*   _2^2 / \varepsilon^2) \\ \mathcal{O}(mn\min\{m, n\})$
Flow polytope Birkhoff polytope Permutahedron	$\mathcal{O}(m+n)$ $\mathcal{O}(n^3)$ $\mathcal{O}(n \ln(n))$	$ \begin{array}{l} \tilde{\mathcal{O}}(m^3n + n^2) \\ \mathcal{O}(n^2 d_z^2 / \varepsilon^2) \\ \mathcal{O}(n \ln(n) + n) \end{array} $

• Can we avoid projections?

The Frank-Wolfe algorithm (Frank & Wolfe, 1956) a.k.a. conditional gradient algorithm (Levitin & Polyak, 1966):



The Frank-Wolfe algorithm (Frank & Wolfe, 1956) a.k.a. conditional gradient algorithm (Levitin & Polyak, 1966):



The Frank-Wolfe algorithm (Frank & Wolfe, 1956) a.k.a. conditional gradient algorithm (Levitin & Polyak, 1966):



The Frank-Wolfe algorithm (Frank & Wolfe, 1956) a.k.a. conditional gradient algorithm (Levitin & Polyak, 1966):



The Frank-Wolfe algorithm (Frank & Wolfe, 1956) a.k.a. conditional gradient algorithm (Levitin & Polyak, 1966):



The Frank-Wolfe algorithm (Frank & Wolfe, 1956) a.k.a. conditional gradient algorithm (Levitin & Polyak, 1966):



The Frank-Wolfe algorithm (Frank & Wolfe, 1956) a.k.a. conditional gradient algorithm (Levitin & Polyak, 1966):

AlgorithmFrank-Wolfe (FW)Input: $x_0 \in C, \ \gamma_t \in [0, 1]$ 1:for t = 0 to T - 1 do2: $v_t \leftarrow \arg\min(v, \nabla f(x_t))$ 3: $x_{t+1} \leftarrow x_t + \gamma_t(v_t - x_t)$ 



•  $x_{t+1}$  is obtained by convex combination of  $x_t \in C$  and  $v_t \in C$ , thus  $x_{t+1} \in C$ 

The Frank-Wolfe algorithm (Frank & Wolfe, 1956) a.k.a. conditional gradient algorithm (Levitin & Polyak, 1966):



- $x_{t+1}$  is obtained by convex combination of  $x_t \in C$  and  $v_t \in C$ , thus  $x_{t+1} \in C$
- FW uses linear minimizations (the "FW oracle") instead of projections

The Frank-Wolfe algorithm (Frank & Wolfe, 1956) a.k.a. conditional gradient algorithm (Levitin & Polyak, 1966):



- $x_{t+1}$  is obtained by convex combination of  $x_t \in C$  and  $v_t \in C$ , thus  $x_{t+1} \in C$
- FW uses linear minimizations (the "FW oracle") instead of projections
- FW = pick a vertex (using gradient information) and move in that direction

The Frank-Wolfe algorithm (Frank & Wolfe, 1956) a.k.a. conditional gradient algorithm (Levitin & Polyak, 1966):



- $x_{t+1}$  is obtained by convex combination of  $x_t \in C$  and  $v_t \in C$ , thus  $x_{t+1} \in C$
- FW uses linear minimizations (the "FW oracle") instead of projections
- FW = pick a vertex (using gradient information) and move in that direction
- Applications: traffic assignment, computer vision, optimal transport, adversarial learning, etc.

• The first strategy considered historically (Frank & Wolfe 1956; Levitin & Polyak, 1966; Demyanov & Rubinov, 1970) is

$$\gamma_t \leftarrow \min\left\{\frac{\langle x_t - \mathbf{v}_t, \nabla f(x_t) \rangle}{L \|x_t - \mathbf{v}_t\|^2}, 1\right\}$$

and is obtained from the smoothness upper bound:

$$\gamma_t = \operatorname*{arg\,min}_{\gamma \in [0,1]} f(x_t) + \gamma \langle v_t - x_t, \nabla f(x_t) \rangle + \frac{L}{2} \gamma^2 \|v_t - x_t\|_2^2$$

• The first strategy considered historically (Frank & Wolfe 1956; Levitin & Polyak, 1966; Demyanov & Rubinov, 1970) is

$$\gamma_t \leftarrow \min\left\{\frac{\langle x_t - v_t, \nabla f(x_t)\rangle}{L\|x_t - v_t\|^2}, 1\right\}$$

and is obtained from the smoothness upper bound:

$$\gamma_t = \operatorname*{arg\,min}_{\gamma \in [0,1]} f(x_t) + \gamma \langle v_t - x_t, \nabla f(x_t) \rangle + \frac{L}{2} \gamma^2 \|v_t - x_t\|_2^2$$

• The first strategy considered historically (Frank & Wolfe 1956; Levitin & Polyak, 1966; Demyanov & Rubinov, 1970) is

$$\gamma_t \leftarrow \min\left\{\frac{\langle x_t - \mathbf{v}_t, \nabla f(\mathbf{x}_t) \rangle}{L \|\mathbf{x}_t - \mathbf{v}_t\|^2}, 1\right\}$$

and is obtained from the smoothness upper bound:

$$\gamma_t = \operatorname*{arg\,min}_{\gamma \in [0,1]} f(x_t) + \gamma \langle v_t - x_t, \nabla f(x_t) \rangle + \frac{L}{2} \gamma^2 \|v_t - x_t\|_2^2$$

• The first strategy considered historically (Frank & Wolfe 1956; Levitin & Polyak, 1966; Demyanov & Rubinov, 1970) is

$$\gamma_t \leftarrow \min\left\{\frac{\langle x_t - v_t, \nabla f(x_t)\rangle}{L\|x_t - v_t\|^2}, 1\right\}$$

and is obtained from the smoothness upper bound:

$$\gamma_t = \operatorname*{arg\,min}_{\gamma \in [0,1]} f(x_t) + \gamma \langle v_t - x_t, \nabla f(x_t) \rangle + \frac{L}{2} \gamma^2 \|v_t - x_t\|_2^2$$

 Later on, Dunn & Harshbarger (1978) proposed open-loop strategies in the form γ<sub>t</sub> ~ 1/t. The one popularized by Jaggi (2013) is

$$\gamma_t \leftarrow \frac{2}{t+2}$$

Theorem (Frank & Wolfe, 1956; Levitin & Polyak, 1966; Jaggi, 2013)

Let  $C \subset \mathbb{R}^n$  be a compact convex set with diameter D and  $f : \mathbb{R}^n \to \mathbb{R}$  be an L-smooth convex function. Then

$$f(x_t) - \min_{\mathcal{C}} f \leqslant \frac{4LD^2}{t+2}$$

Theorem (Frank & Wolfe, 1956; Levitin & Polyak, 1966; Jaggi, 2013)

Let  $C \subset \mathbb{R}^n$  be a compact convex set with diameter D and  $f : \mathbb{R}^n \to \mathbb{R}$  be an L-smooth convex function. Then

$$f(x_t) - \min_{\mathcal{C}} f \leqslant \frac{4LD^2}{t+2}$$

• The convergence rate cannot be improved (Canon & Cullum, 1968; Jaggi, 2013; Lan, 2013)

Theorem (Frank & Wolfe, 1956; Levitin & Polyak, 1966; Jaggi, 2013)

Let  $C \subset \mathbb{R}^n$  be a compact convex set with diameter D and  $f : \mathbb{R}^n \to \mathbb{R}$  be an L-smooth convex function. Then

$$f(x_t) - \min_{\mathcal{C}} f \leqslant \frac{4LD^2}{t+2}$$

- The convergence rate cannot be improved (Canon & Cullum, 1968; Jaggi, 2013; Lan, 2013)
- The iterates can follow a naive and slow trajectory

Consider the simple problem

$$\min \frac{1}{2} \|x\|_2^2$$
  
s.t.  $x \in \operatorname{conv} \left\{ \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right\}$ 

and 
$$x^* = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$



Consider the simple problem

$$\min \frac{1}{2} \|x\|_2^2$$
  
s.t.  $x \in \operatorname{conv} \left\{ \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right\}$ 

and 
$$x^* = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

• Let 
$$x_0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$



Consider the simple problem

$$\min \frac{1}{2} \|x\|_2^2$$
  
s.t.  $x \in \operatorname{conv} \left\{ \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right\}$ 

and 
$$x^* = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

• Let 
$$x_0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

• FW tries to reach x\* by moving towards vertices



Consider the simple problem

$$\min \frac{1}{2} \|x\|_2^2$$
  
s.t.  $x \in \operatorname{conv} \left\{ \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right\}$ 

and 
$$x^* = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

• Let 
$$x_0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

• FW tries to reach x\* by moving towards vertices


Consider the simple problem

$$\min \frac{1}{2} \|x\|_2^2$$
  
s.t.  $x \in \operatorname{conv} \left\{ \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right\}$ 

and 
$$x^* = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

• Let 
$$x_0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$



Consider the simple problem

$$\begin{split} &\min \, \frac{1}{2} \|x\|_2^2 \\ &\text{s.t. } x \in \operatorname{conv} \left\{ \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right\} \end{split}$$

and 
$$x^* = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

• Let 
$$x_0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$



Consider the simple problem

$$\min \frac{1}{2} \|x\|_2^2$$
  
s.t.  $x \in \operatorname{conv} \left\{ \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right\}$ 

and 
$$x^* = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

• Let 
$$x_0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$



Consider the simple problem

$$\min \frac{1}{2} \|x\|_2^2$$
  
s.t.  $x \in \operatorname{conv} \left\{ \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right\}$ 

and 
$$x^* = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

• Let 
$$x_0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$



Consider the simple problem

$$\min \frac{1}{2} \|x\|_2^2$$
  
s.t.  $x \in \operatorname{conv} \left\{ \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right\}$ 

and 
$$x^* = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

• Let 
$$x_0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$



Consider the simple problem

$$\begin{split} &\min \, \frac{1}{2} \|x\|_2^2 \\ &\text{s.t. } x \in \operatorname{conv} \left\{ \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right\} \end{split}$$

and 
$$x^* = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

• Let 
$$x_0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$



Consider the simple problem

$$\begin{split} &\min \, \frac{1}{2} \|x\|_2^2 \\ &\text{s.t. } x \in \operatorname{conv} \left\{ \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right\} \end{split}$$

and 
$$x^* = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

• Let 
$$x_0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$



Consider the simple problem

$$\begin{split} &\min \, \frac{1}{2} \|x\|_2^2 \\ &\text{s.t. } x \in \operatorname{conv} \left\{ \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right\} \end{split}$$

and 
$$x^* = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

• Let 
$$x_0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

- FW tries to reach x\* by moving towards vertices
- This yields an inefficient zig-zagging trajectory



 Away-Step Frank-Wolfe (AFW) (Wolfe, 1970; Guélat & Marcotte, 1986; Lacoste-Julien & Jaggi, 2015): enhances FW by allowing to also move away from vertices

 Away-Step Frank-Wolfe (AFW) (Wolfe, 1970; Guélat & Marcotte, 1986; Lacoste-Julien & Jaggi, 2015): enhances FW by allowing to also move away from vertices



 Away-Step Frank-Wolfe (AFW) (Wolfe, 1970; Guélat & Marcotte, 1986; Lacoste-Julien & Jaggi, 2015): enhances FW by allowing to also move away from vertices



 Decomposition-Invariant Pairwise Conditional Gradient (DICG) (Garber & Meshi, 2016): memory-free variant of AFW

 Away-Step Frank-Wolfe (AFW) (Wolfe, 1970; Guélat & Marcotte, 1986; Lacoste-Julien & Jaggi, 2015): enhances FW by allowing to also move away from vertices



- Decomposition-Invariant Pairwise Conditional Gradient (DICG) (Garber & Meshi, 2016): memory-free variant of AFW
- Blended Conditional Gradients (BCG) (Braun et al., 2019): blends FCFW and FW

• Can we speed up FW in a simple way?

- Can we speed up FW in a simple way?
- Rule of thumb in optimization: follow the steepest direction

- Can we speed up FW in a simple way?
- Rule of thumb in optimization: follow the steepest direction

Idea:

• Speed up FW by moving in a direction better aligned with  $-\nabla f(x_t)$ 

- Can we speed up FW in a simple way?
- Rule of thumb in optimization: follow the steepest direction

Idea:

- Speed up FW by moving in a direction better aligned with  $-\nabla f(x_t)$
- Build this direction by using  ${\mathcal C}$  to maintain the projection-free property



• 
$$v_0 \in \arg \max_{v \in C} \langle v, -\nabla f(x_t) \rangle$$
  
 $\lambda_0 u_0 = \frac{\langle v_0 - x_t, -\nabla f(x_t) \rangle}{\|v_0 - x_t\|^2} (v_0 - x_t)$   
 $r_1 = -\nabla f(x_t) - \lambda_0 u_0$ 



• 
$$v_0 \in \arg \max_{v \in C} \langle v, -\nabla f(x_t) \rangle$$
  
 $\lambda_0 u_0 = \frac{\langle v_0 - x_t, -\nabla f(x_t) \rangle}{\|v_0 - x_t\|^2} (v_0 - x_t)$   
 $r_1 = -\nabla f(x_t) - \lambda_0 u_0$ 



 How can we build a direction better aligned with −∇f(x<sub>t</sub>) and that allows to update x<sub>t+1</sub> without projection?

• 
$$v_0 \in \arg \max_{v \in C} \langle v, -\nabla f(x_t) \rangle$$
  
 $\lambda_0 u_0 = \frac{\langle v_0 - x_t, -\nabla f(x_t) \rangle}{\|v_0 - x_t\|^2} (v_0 - x_t)$   
 $r_1 = -\nabla f(x_t) - \lambda_0 u_0$ 

•  $v_1 \in \arg \max_{v \in C} \langle v, r_1 \rangle$   $\lambda_1 u_1 = \frac{\langle v_1 - x_t, r_1 \rangle}{\|v_1 - x_t\|^2} (v_1 - x_t)$  $r_2 = r_1 - \lambda_1 u_1$ 



• 
$$v_0 \in \arg \max_{v \in C} \langle v, -\nabla f(x_t) \rangle$$
  
 $\lambda_0 u_0 = \frac{\langle v_0 - x_t, -\nabla f(x_t) \rangle}{\|v_0 - x_t\|^2} (v_0 - x_t)$   
 $r_1 = -\nabla f(x_t) - \lambda_0 u_0$ 

- $v_1 \in \arg \max_{v \in C} \langle v, r_1 \rangle$   $\lambda_1 u_1 = \frac{\langle v_1 - x_t, r_1 \rangle}{\|v_1 - x_t\|^2} (v_1 - x_t)$  $r_2 = r_1 - \lambda_1 u_1$
- We could continue:  $v_2 \in \arg \max_{v \in C} \langle v, r_2 \rangle$



• 
$$v_0 \in \arg \max_{v \in C} \langle v, -\nabla f(x_t) \rangle$$
  
 $\lambda_0 u_0 = \frac{\langle v_0 - x_t, -\nabla f(x_t) \rangle}{\|v_0 - x_t\|^2} (v_0 - x_t)$   
 $r_1 = -\nabla f(x_t) - \lambda_0 u_0$ 

- $v_1 \in \arg \max_{v \in C} \langle v, r_1 \rangle$   $\lambda_1 u_1 = \frac{\langle v_1 - x_t, r_1 \rangle}{\|v_1 - x_t\|^2} (v_1 - x_t)$  $r_2 = r_1 - \lambda_1 u_1$
- We could continue:  $v_2 \in \arg \max_{v \in C} \langle v, r_2 \rangle$
- $d = \lambda_0 u_0 + \lambda_1 u_1$



• 
$$v_0 \in \arg \max_{v \in C} \langle v, -\nabla f(x_t) \rangle$$
  
 $\lambda_0 u_0 = \frac{\langle v_0 - x_t, -\nabla f(x_t) \rangle}{\|v_0 - x_t\|^2} (v_0 - x_t)$   
 $r_1 = -\nabla f(x_t) - \lambda_0 u_0$ 

- $v_1 \in \arg \max_{v \in C} \langle v, r_1 \rangle$   $\lambda_1 u_1 = \frac{\langle v_1 - x_t, r_1 \rangle}{\|v_1 - x_t\|^2} (v_1 - x_t)$  $r_2 = r_1 - \lambda_1 u_1$
- We could continue:  $v_2 \in \arg \max_{v \in C} \langle v, r_2 \rangle$
- $d = \lambda_0 u_0 + \lambda_1 u_1$
- $g_t = d/(\lambda_0 + \lambda_1)$



• 
$$v_0 \in \arg \max_{v \in C} \langle v, -\nabla f(x_t) \rangle$$
  
 $\lambda_0 u_0 = \frac{\langle v_0 - x_t, -\nabla f(x_t) \rangle}{\|v_0 - x_t\|^2} (v_0 - x_t)$   
 $r_1 = -\nabla f(x_t) - \lambda_0 u_0$ 

- $v_1 \in \arg \max_{v \in C} \langle v, r_1 \rangle$   $\lambda_1 u_1 = \frac{\langle v_1 - x_t, r_1 \rangle}{\|v_1 - x_t\|^2} (v_1 - x_t)$  $r_2 = r_1 - \lambda_1 u_1$
- We could continue:  $v_2 \in \arg \max_{v \in C} \langle v, r_2 \rangle$
- $d = \lambda_0 u_0 + \lambda_1 u_1$
- $g_t = d/(\lambda_0 + \lambda_1)$
- The boosted direction  $g_t$  is better aligned with  $-\nabla f(x_t)$  than is the FW direction  $v_0 x_t$



• 
$$v_0 \in \arg \max_{v \in C} \langle v, -\nabla f(x_t) \rangle$$
  
 $\lambda_0 u_0 = \frac{\langle v_0 - x_t, -\nabla f(x_t) \rangle}{\|v_0 - x_t\|^2} (v_0 - x_t)$   
 $r_1 = -\nabla f(x_t) - \lambda_0 u_0$ 

- $v_1 \in \arg \max_{v \in C} \langle v, r_1 \rangle$   $\lambda_1 u_1 = \frac{\langle v_1 - x_t, r_1 \rangle}{\|v_1 - x_t\|^2} (v_1 - x_t)$  $r_2 = r_1 - \lambda_1 u_1$
- We could continue:  $v_2 \in \arg \max_{v \in C} \langle v, r_2 \rangle$
- $d = \lambda_0 u_0 + \lambda_1 u_1$
- $g_t = d/(\lambda_0 + \lambda_1)$



$$x_{t+1} = x_t + \gamma_t g_t$$
 for all  $\gamma_t \in [0, 1]$ 



## A generic boosting procedure

#### **Algorithm** Boosting procedure Boost( $\mathbf{d}, \mathbf{z}, \mathcal{K}, \delta$ )

Input: 
$$d \neq 0, z \in C, K \in \mathbb{N} \setminus \{0\}, \delta \in ]0, 1[$$

 1:  $d_0 \leftarrow 0, \Lambda \leftarrow 0$ 

 2: for  $k = 0$  to  $K - 1$  do

 3:  $r_k \leftarrow d - d_k$ 
 > k-th residual

 4:  $v_k \leftarrow \arg \max_{v \in C} \langle v, r_k \rangle$ 
 > FW oracle

 5:  $u_k \leftarrow v_k - z$ 
 > FW oracle

 6:  $\lambda_k \leftarrow \langle u_k, r_k \rangle / || u_k ||_2^2$ 
 >

 7:  $d'_k \leftarrow d_k + \lambda_k u_k$ 
 >  $\delta$  then

 9:  $d_{k+1} \leftarrow d'_k$ 
 >  $\delta$  then

 10:  $\Lambda \leftarrow \Lambda + \lambda_k$ 
 >  $\delta$  then

 11: else
 > exit k-loop

 13:  $g \leftarrow d_k / \Lambda$ 
 > normalization

• 
$$\cos(\hat{d}, \mathbf{d}) = \frac{\langle \hat{d}, \mathbf{d} \rangle}{\|\hat{d}\|_2 \|\mathbf{d}\|_2}$$
 if  $\hat{d} \neq 0$ , else  $-1$  if  $\hat{d} = 0$ 

## A generic boosting procedure

#### **Algorithm** Boosting procedure $Boost(\mathbf{d}, \mathbf{z}, K, \delta)$

Input: 
$$d \neq 0, z \in C, K \in \mathbb{N} \setminus \{0\}, \delta \in ]0, 1[$$

 1:  $d_0 \leftarrow 0, \Lambda \leftarrow 0$ 

 2: for  $k = 0$  to  $K - 1$  do

 3:  $r_k \leftarrow d - d_k$ 
 $\triangleright$  k-th residual

 4:  $v_k \leftarrow \arg \max_{v \in C} \langle v, r_k \rangle$ 
 $\triangleright$  FW oracle

 5:  $u_k \leftarrow v_k - z$ 
 $\circ$ 

 6:  $\lambda_k \leftarrow \langle u_k, r_k \rangle / || u_k ||_2^2$ 
 $\rho$ 

 7:  $d'_k \leftarrow d_k + \lambda_k u_k$ 
 $\circ$ 

 8: if  $\cos(d'_k, d) - \cos(d_k, d) \ge \delta$  then
  $\circ$ 

 9:  $d_{k+1} \leftarrow d'_k$ 
 $\circ$ 

 10:  $\Lambda \leftarrow \Lambda + \lambda_k$ 
 $\circ$ 

 11: else
  $\circ$  exit k-loop

 13:  $g \leftarrow d_k / \Lambda$ 
 $\triangleright$  normalization

• 
$$\cos(\hat{d}, \mathbf{d}) = \frac{\langle \hat{d}, \mathbf{d} \rangle}{\|\hat{d}\|_2 \|\mathbf{d}\|_2}$$
 if  $\hat{d} \neq 0$ , else  $-1$  if  $\hat{d} = 0$ 

• The stopping criterion is an alignment improvement condition (typically  $\delta \leftarrow 10^{-3}$  and  $K \leftarrow +\infty$ )

## A generic boosting procedure

#### **Algorithm** Boosting procedure Boost( $\mathbf{d}, \mathbf{z}, K, \delta$ ) **Input:** $\mathbf{d} \neq \mathbf{0}, \mathbf{z} \in \mathcal{C}, K \in \mathbb{N} \setminus \{\mathbf{0}\}, \delta \in [0, 1[$ 1: $d_0 \leftarrow 0, \Lambda \leftarrow 0$ 2: for k = 0 to K - 1 do 3: $r_k \leftarrow \mathbf{d} - d_k$ $\triangleright$ k-th residual 4: $v_k \leftarrow \arg \max_{v \in C} \langle v, r_k \rangle$ ▷ FW oracle 5: $u_k \leftarrow v_k - \mathbf{z}$ 6: $\lambda_k \leftarrow \langle u_k, r_k \rangle / \|u_k\|_2^2$ 7: $d'_k \leftarrow d_k + \lambda_k u_k$ if $\cos(d'_k, \mathbf{d}) - \cos(d_k, \mathbf{d}) \ge \delta$ then 8: $d_{k+1} \leftarrow d'_k$ 9: $\Lambda \leftarrow \Lambda + \lambda_{\mu}$ 10: 11: else break 12: $\triangleright$ exit k-loop 13: $g \leftarrow d_k / \Lambda$ ▷ normalization

• 
$$\cos(\hat{d}, \mathbf{d}) = \frac{\langle \hat{d}, \mathbf{d} \rangle}{\|\hat{d}\|_2 \|\mathbf{d}\|_2}$$
 if  $\hat{d} \neq 0$ , else  $-1$  if  $\hat{d} = 0$ 

• The stopping criterion is an alignment improvement condition (typically  $\delta \leftarrow 10^{-3}$  and  $K \leftarrow +\infty$ )

AlgorithmFrank-Wolfe (FW)Input: $x_0 \in C, \ \gamma_t \in [0, 1]$ 1:for t = 0 to T - 1 do2: $v_t \leftarrow \arg\min \langle \nabla f(x_t), v \rangle$ 3: $x_{t+1} \leftarrow x_t + \gamma_t (v_t - x_t)$ 

**Algorithm** Boosted Frank-Wolfe (BoostFW) **Input:**  $x_0 \in C$ ,  $\gamma_t \in [0, 1]$ ,  $K \in \mathbb{N} \setminus \{0\}$ ,  $\delta \in [0, 1]$ 

1: for t = 0 to T - 1 do

2: 
$$g_t \leftarrow \text{Boost}(-\nabla f(x_t), x_t, K, \delta)$$

3:  $x_{t+1} \leftarrow x_t + \gamma_t g_t$ 

AlgorithmFrank-Wolfe (FW)Input: $x_0 \in C, \ \gamma_t \in [0, 1]$ 1:for t = 0 to T - 1 do2: $v_t \leftarrow \underset{v \in C}{\operatorname{arg\,min}} \langle \nabla f(x_t), v \rangle$ 3: $x_{t+1} \leftarrow x_t + \gamma_t (v_t - x_t)$ 

Algorithm Boosted Frank-Wolfe (BoostFW)

Input:  $x_0 \in C$ ,  $\gamma_t \in [0, 1]$ ,  $K \in \mathbb{N} \setminus \{0\}$ ,  $\delta \in ]0, 1[$ 1: for t = 0 to T - 1 do

- 2:  $g_t \leftarrow \text{Boost}(-\nabla f(x_t), x_t, K, \delta)$
- 3:  $x_{t+1} \leftarrow x_t + \gamma_t g_t$

AlgorithmFrank-Wolfe (FW)Input: $x_0 \in C, \ \gamma_t \in [0, 1]$ 1:for t = 0 to T - 1 do2: $v_t \leftarrow \arg\min(\nabla f(x_t), v)$ 3: $x_{t+1} \leftarrow x_t + \gamma_t(v_t - x_t)$ 

Algorithm Boosted Frank-Wolfe (BoostFW)

Input:  $x_0 \in C$ ,  $\gamma_t \in [0, 1]$ ,  $K \in \mathbb{N} \setminus \{0\}$ ,  $\delta \in ]0, 1[$ 1: for t = 0 to T - 1 do

2: 
$$g_t \leftarrow \text{Boost}(-\nabla f(x_t), x_t, K, \delta)$$

3:  $x_{t+1} \leftarrow x_t + \gamma_t g_t$ 



x<sub>0</sub> x<sub>1</sub> x<sub>3</sub> x<sub>2</sub> x<sup>4</sup>

Algorithm Boosted Frank-Wolfe (BoostFW) Input:  $x_0 \in C$ ,  $\gamma_t \in [0, 1]$ ,  $K \in \mathbb{N} \setminus \{0\}$ ,  $\delta \in ]0, 1[$ 1: for t = 0 to T - 1 do

2: 
$$g_t \leftarrow \text{Boost}(-\nabla f(x_t), x_t, K, \delta)$$

3: 
$$x_{t+1} \leftarrow x_t + \gamma_t g_t$$









• What is the convergence rate of BoostFW?



AlgorithmBoosted Frank-Wolfe (BoostFW)Input: $x_0 \in C, \gamma_t \in [0, 1], K \in \mathbb{N} \setminus \{0\}, \delta \in ]0, 1[$ 1:for t = 0 to T - 1 do2: $g_t \leftarrow \text{Boost}(-\nabla f(x_t), x_t, K, \delta)$ 3: $x_{t+1} \leftarrow x_t + \gamma_t g_t$ 



 $x^* = x_1$ 

- What is the convergence rate of BoostFW?
- Is BoostFW expensive in practice?



AlgorithmBoosted Frank-Wolfe (BoostFW)Input: $x_0 \in C, \ \gamma_t \in [0, 1], \ K \in \mathbb{N} \setminus \{0\}, \ \delta \in ]0, 1[$ 1:for t = 0 to T - 1 do2: $g_t \leftarrow \text{Boost}(-\nabla f(x_t), x_t, K, \delta)$ 3: $x_{t+1} \leftarrow x_t + \gamma_t g_t$ 



$$x^* = x_1$$

- What is the convergence rate of BoostFW?
- Is BoostFW expensive in practice?
- How does it compare to the state of the art?
• Let  $N_t$  be the number of iterations up to t for which at least 2 rounds of alignment were performed (FW = always 1 round) with a step-size < 1

• Let N<sub>t</sub> be the number of iterations up to t for which at least 2 rounds of alignment were performed (FW = always 1 round) with a step-size < 1

#### Theorem

Let  $C \subset \mathbb{R}^n$  be a compact convex set with diameter D and  $f : \mathbb{R}^n \to \mathbb{R}$  be an L-smooth, convex, and  $\mu$ -gradient dominated function, and let  $x_0 \leftarrow \arg\min_{v \in C} \langle v, \nabla f(y) \rangle$  for some  $y \in C$  and  $\gamma_t \leftarrow \min\left\{\frac{\langle g_t, -\nabla f(x_t) \rangle}{L \|g_t\|_2^2}, 1\right\}$ . Suppose that  $N_t \ge \omega t$  where  $\omega > 0$ . Then

$$f(x_t) - \min_{\mathcal{C}} f \leqslant \frac{LD^2}{2} \exp\left(-\delta^2 \frac{\mu}{L} \omega t\right)$$

• Let  $N_t$  be the number of iterations up to t for which at least 2 rounds of alignment were performed (FW = always 1 round) with a step-size < 1

#### Theorem

Let  $C \subset \mathbb{R}^n$  be a compact convex set with diameter D and  $f : \mathbb{R}^n \to \mathbb{R}$  be an L-smooth, convex, and  $\mu$ -gradient dominated function, and let  $x_0 \leftarrow \arg\min_{v \in C} \langle v, \nabla f(y) \rangle$  for some  $y \in C$  and  $\gamma_t \leftarrow \min\left\{\frac{\langle g_t, -\nabla f(x_t) \rangle}{L \|g_t\|_2^2}, 1\right\}$ . Suppose that  $N_t \ge \omega t$  where  $\omega > 0$ . Then

$$f(x_t) - \min_{\mathcal{C}} f \leqslant \frac{LD^2}{2} \exp\left(-\delta^2 \frac{\mu}{L} \omega t\right)$$

 The assumption N<sub>t</sub> ≥ ωt simply states that N<sub>t</sub> is nonnegligeable, i.e., that the boosting procedure is active

• Let N<sub>t</sub> be the number of iterations up to t for which at least 2 rounds of alignment were performed (FW = always 1 round) with a step-size < 1

#### Theorem

Let  $C \subset \mathbb{R}^n$  be a compact convex set with diameter D and  $f : \mathbb{R}^n \to \mathbb{R}$  be an L-smooth, convex, and  $\mu$ -gradient dominated function, and let  $x_0 \leftarrow \arg\min_{v \in C} \langle v, \nabla f(y) \rangle$  for some  $y \in C$  and  $\gamma_t \leftarrow \min\left\{\frac{\langle g_t, -\nabla f(x_t) \rangle}{L \|g_t\|_2^2}, 1\right\}$ . Suppose that  $N_t \ge \omega t$  where  $\omega > 0$ . Then

$$f(x_t) - \min_{\mathcal{C}} f \leqslant \frac{LD^2}{2} \exp\left(-\delta^2 \frac{\mu}{L} \omega t\right)$$

- The assumption N<sub>t</sub> ≥ ωt simply states that N<sub>t</sub> is nonnegligeable, i.e., that the boosting procedure is active
- Else, BoostFW reduces to FW and the convergence rate is  $\frac{4LD^2}{t+2}$

• Let  $N_t$  be the number of iterations up to t for which at least 2 rounds of alignment were performed (FW = always 1 round) with a step-size < 1

#### Theorem

Let  $C \subset \mathbb{R}^n$  be a compact convex set with diameter D and  $f : \mathbb{R}^n \to \mathbb{R}$  be an L-smooth, convex, and  $\mu$ -gradient dominated function, and let  $x_0 \leftarrow \arg\min_{v \in C} \langle v, \nabla f(y) \rangle$  for some  $y \in C$  and  $\gamma_t \leftarrow \min\left\{\frac{\langle g_t, -\nabla f(x_t) \rangle}{L \|g_t\|_2^2}, 1\right\}$ . Suppose that  $N_t \ge \omega t$  where  $\omega > 0$ . Then

$$f(x_t) - \min_{\mathcal{C}} f \leqslant \frac{LD^2}{2} \exp\left(-\delta^2 \frac{\mu}{L} \omega t\right)$$

- The assumption N<sub>t</sub> ≥ ωt simply states that N<sub>t</sub> is nonnegligeable, i.e., that the boosting procedure is active
- Else, BoostFW reduces to FW and the convergence rate is  $\frac{4LD^2}{t+2}$
- In practice,  $N_t pprox t$  (so  $\omega \lesssim 1$ )

 We compare BoostFW to AFW, BCG, and DICG on a series of experiments involving various objective functions and feasible regions

 We compare BoostFW to AFW, BCG, and DICG on a series of experiments involving various objective functions and feasible regions

$$\min_{x \in \mathbb{R}^n} \frac{1}{m} \sum_{i=1}^m \ln(1 + \exp(-y_i \langle a_i, x \rangle))$$
s.t.  $\|x\|_1 \leqslant \tau$ 

$$\begin{split} \min_{\substack{X \in \mathbb{R}^{m \times n}}} \frac{1}{|\mathcal{I}|} \sum_{(i,j) \in \mathcal{I}} h_{\rho}(Y_{i,j} - X_{i,j}) \\ \text{s.t. } \|X\|_{\text{nuc}} \leqslant \tau \end{split}$$

 We compare BoostFW to AFW, BCG, and DICG on a series of experiments involving various objective functions and feasible regions

$$\begin{array}{l} \min_{x \in \mathbb{R}^{n}} \sum_{a \in \mathcal{A}} \tau_{a} x_{a} \left(1 + 0.03 \left(\frac{x_{a}}{c_{a}}\right)^{4}\right) \\ \text{s.t. } \|y\|_{1} \leqslant \tau \end{array} \\ \begin{array}{l} \text{s.t. } \|x\|_{1} \leqslant \tau \end{array} \\ \begin{array}{l} \text{s.t. } x_{a} = \sum_{r \in \mathcal{R}} \mathbb{1}_{\{a \in r\}} y_{r} \quad a \in \mathcal{A} \\ \sum_{r \in \mathcal{R}_{i,j}} y_{r} = d_{i,j} \quad (i,j) \in \mathcal{S} \\ y_{r} \geqslant 0 \qquad r \in \mathcal{R}_{i,j}, (i,j) \in \mathcal{S} \end{array}$$

$$\min_{\boldsymbol{x}\in\mathbb{R}^n} \frac{1}{m} \sum_{i=1}^m \ln(1 + \exp(-y_i \langle \boldsymbol{a}_i, \boldsymbol{x} \rangle))$$
s.t.  $\|\boldsymbol{x}\|_1 \leq \tau$ 

$$\min_{\boldsymbol{X}\in\mathbb{R}^{m \times n}} \frac{1}{|\mathcal{I}|} \sum_{(i,j)\in\mathcal{I}} h_{\rho}(\boldsymbol{Y}_{i,j} - \boldsymbol{X}_{i,j})$$
s.t.  $\|\boldsymbol{X}\|_{\text{nuc}} \leq \tau$ 

 For BoostFW and AFW we also run the line search-free strategies and label them with an "L"



Sparse signal recovery





• Sparse logistic regression on the Gisette dataset

40

80



• Collaborative filtering on the MovieLens 100k dataset



• DICG is known to perform particularly well on the video co-localization experiment (YouTube-Objects dataset)

- DICG is known to perform particularly well on the video co-localization experiment (YouTube-Objects dataset)
- BoostDICG: application of our method to DICG



- DICG is known to perform particularly well on the video co-localization experiment (YouTube-Objects dataset)
- BoostDICG: application of our method to DICG



• (details)

DICG

$$a_t \leftarrow \text{away vertex}$$
$$v_t \leftarrow \underset{v \in \mathcal{C}}{\arg\min} \langle v, \nabla f(x_t) \rangle$$
$$x_{t+1} \leftarrow x_t + \gamma_t (v_t - a_t)$$

- DICG is known to perform particularly well on the video co-localization experiment (YouTube-Objects dataset)
- BoostDICG: application of our method to DICG



• (details)

DICG

$$a_t \leftarrow \text{away vertex}$$

$$v_t \leftarrow \underset{v \in C}{\operatorname{arg min}} \langle v, \nabla f(x_t) \rangle$$

$$x_{t+1} \leftarrow x_t + \gamma_t(v_t - a_t)$$

### BoostDICG

$$a_t \leftarrow \text{away vertex}$$
  
 $g_t \leftarrow \text{Boost}(-\nabla f(x_t), a_t, K, \delta)$   
 $x_{t+1} \leftarrow x_t + \gamma_t g_t$ 

- DICG is known to perform particularly well on the video co-localization experiment (YouTube-Objects dataset)
- BoostDICG: application of our method to DICG



• (details)

DICG

$$a_t \leftarrow \text{away vertex}$$

$$v_t \leftarrow \underset{v \in C}{\operatorname{arg min}} \langle v, \nabla f(x_t) \rangle$$

$$x_{t+1} \leftarrow x_t + \gamma_t(v_t - a_t)$$

### BoostDICG

$$a_t \leftarrow \text{away vertex}$$
  
 $g_t \leftarrow \text{Boost}(-\nabla f(x_t), a_t, K, \delta)$   
 $x_{t+1} \leftarrow x_t + \gamma_t g_t$ 

- DICG is known to perform particularly well on the video co-localization experiment (YouTube-Objects dataset)
- BoostDICG: application of our method to DICG



• (details)

DICG

$$a_t \leftarrow \text{away vertex}$$

$$v_t \leftarrow \underset{v \in C}{\operatorname{arg min}} \langle v, \nabla f(x_t) \rangle$$

$$x_{t+1} \leftarrow x_t + \gamma_t(v_t - a_t)$$

### BoostDICG

$$a_t \leftarrow \text{away vertex}$$
  
 $g_t \leftarrow \text{Boost}(-\nabla f(x_t), a_t, K, \delta)$   
 $x_{t+1} \leftarrow x_t + \gamma_t g_t$ 

### Large-scale optimization

Consider

$$\min \left\{ f(x) \coloneqq \frac{1}{m} \sum_{i=1}^{m} f_i(x) \right\}$$
  
s.t.  $x \in C$ 

where

- $\mathcal{C} \subset \mathbb{R}^n$  is a compact convex set
- $f_1, \ldots, f_m \colon \mathbb{R}^n \to \mathbb{R}$  are smooth (non)convex functions
- *m* ≫ 1 is very large

### Large-scale optimization

Consider

$$\min \left\{ f(x) \coloneqq \frac{1}{m} \sum_{i=1}^{m} f_i(x) \right\}$$
  
s.t.  $x \in C$ 

where

- $\mathcal{C} \subset \mathbb{R}^n$  is a compact convex set
- $f_1, \ldots, f_m \colon \mathbb{R}^n \to \mathbb{R}$  are smooth (non)convex functions
- *m* ≫ 1 is very large

Computing f(x) or  $\nabla f(x)$  is too expensive

- Cannot use line search
- More efficient to use an estimator  $\tilde{\nabla} f(x)$  to get approximate (but cheap) gradient information

# TemplateStochastic Frank-WolfeInput: $x_0 \in C$ , $\gamma_t \in [0, 1]$ 1:for t = 0 to T - 1 do2:Update the gradient estimator $\tilde{\nabla}f(x_t)$ 3: $v_t \leftarrow \arg\min_{v \in C} \langle v, \tilde{\nabla}f(x_t) \rangle$ 4: $x_{t+1} \leftarrow x_t + \gamma_t (v_t - x_t)$



• Use a stochastic estimator  $\tilde{\nabla} f(x_t)$ 

# TemplateStochastic Frank-WolfeInput: $x_0 \in C, \ \gamma_t \in [0, 1]$ 1: for t = 0 to T - 1 do2: Update the gradient estimator $\tilde{\nabla}f(x_t)$ 3: $v_t \leftarrow \arg\min_{v \in C} \langle v, \tilde{\nabla}f(x_t) \rangle$ 4: $x_{t+1} \leftarrow x_t + \gamma_t(v_t - x_t)$

- Use a stochastic estimator  $\tilde{\nabla} f(x_t)$
- The *vanilla* Stochastic Frank-Wolfe algorithm (SFW) estimates the gradient by averaging over a minibatch of size *b<sub>t</sub>*:

$$ilde{
abla} f(x_t) \leftarrow rac{1}{b_t} \sum_{j=1}^{b_t} 
abla f_{i_j}(x_t) \quad ext{where} \quad i_1, \dots, i_{b_t} \overset{ ext{i.i.d.}}{\sim} \mathcal{U}(\llbracket 1, m \rrbracket)$$

# TemplateStochastic Frank-WolfeInput: $x_0 \in C, \ \gamma_t \in [0, 1]$ 1:for t = 0 to T - 1 do2:Update the gradient estimator $\tilde{\nabla}f(x_t)$ 3: $v_t \leftarrow \arg\min_{v \in C} \langle v, \tilde{\nabla}f(x_t) \rangle$ 4: $x_{t+1} \leftarrow x_t + \gamma_t(v_t - x_t)$

- Use a stochastic estimator  $\tilde{\nabla} f(x_t)$
- The *vanilla* Stochastic Frank-Wolfe algorithm (SFW) estimates the gradient by averaging over a minibatch of size *b<sub>t</sub>*:

$$ilde{
abla} f(x_t) \leftarrow rac{1}{b_t} \sum_{j=1}^{b_t} 
abla f_{i_j}(x_t) \quad ext{where} \quad i_1, \dots, i_{b_t} \overset{ ext{i.i.d.}}{\sim} \mathcal{U}(\llbracket 1, m \rrbracket)$$

• SFW converges with rate  $\mathcal{O}(1/t)$  when  $b_t \sim t^2$  (Hazan & Luo, 2016)

# TemplateStochastic Frank-WolfeInput: $x_0 \in C, \ \gamma_t \in [0, 1]$ 1: for t = 0 to T - 1 do2: Update the gradient estimator $\tilde{\nabla}f(x_t)$ 3: $v_t \leftarrow \arg\min_{v \in C} \langle v, \tilde{\nabla}f(x_t) \rangle$ 4: $x_{t+1} \leftarrow x_t + \gamma_t(v_t - x_t)$

- Use a stochastic estimator  $\tilde{\nabla} f(x_t)$
- The *vanilla* Stochastic Frank-Wolfe algorithm (SFW) estimates the gradient by averaging over a minibatch of size *b<sub>t</sub>*:

$$\tilde{\nabla} f(\mathbf{x}_t) \leftarrow \frac{1}{b_t} \sum_{j=1}^{b_t} \nabla f_{i_j}(\mathbf{x}_t) \quad \text{where} \quad i_1, \dots, i_{b_t} \overset{\text{i.i.d.}}{\sim} \mathcal{U}(\llbracket 1, m \rrbracket)$$

- SFW converges with rate  $\mathcal{O}(1/t)$  when  $b_t \sim t^2$  (Hazan & Luo, 2016)
- New variants require only  $b_t \sim t$  or  $b_t \sim 1$

Simultaneously proposed by Duchi et al. (2011) and McMahan & Streeter (2010):

AlgorithmAdaptive Gradient (AdaGrad)Input:  $x_0 \in C, \ \delta > 0, \ \eta > 0$ 1: for t = 0 to T - 1 do2: Update the gradient estimator  $\tilde{\nabla}f(x_t)$ 3:  $H_t \leftarrow \text{diag}\left(\delta 1 + \sqrt{\sum_{s=0}^t \tilde{\nabla}f(x_s)^2}\right)$ 4:  $x_{t+1} \leftarrow \arg\min_{x \in C} \eta \langle x, \tilde{\nabla}f(x_t) \rangle + \frac{1}{2} ||x - x_t||_{H_t}^2$ 

Simultaneously proposed by Duchi et al. (2011) and McMahan & Streeter (2010):

 Algorithm
 Adaptive Gradient (AdaGrad)

 Input:
  $x_0 \in C, \ \delta > 0, \ \eta > 0$  

 1:
 for t = 0 to T - 1 do

 2:
 Update the gradient estimator  $\tilde{\nabla}f(x_t)$  

 3:
  $H_t \leftarrow \text{diag}\left(\delta 1 + \sqrt{\sum_{s=0}^t \tilde{\nabla}f(x_s)^2}\right)$  

 4:
  $x_{t+1} \leftarrow \underset{x \in C}{\arg\min \eta \langle x, \tilde{\nabla}f(x_t) \rangle + \frac{1}{2} ||x - x_t||^2_{H_t}$ 

Simultaneously proposed by Duchi et al. (2011) and McMahan & Streeter (2010):

 $\begin{array}{ll} \hline \textbf{Algorithm} & \text{Adaptive Gradient (AdaGrad)} \\ \hline \textbf{Input:} & x_0 \in \mathcal{C}, \ \delta > 0, \ \eta > 0 \\ 1: \ \textbf{for} \ t = 0 \ \textbf{to} \ \mathcal{T} - 1 \ \textbf{do} \\ 2: & \text{Update the gradient estimator } \tilde{\nabla}f(x_t) \\ 3: & H_t \leftarrow \text{diag} \left( \delta 1 + \sqrt{\sum_{s=0}^t \tilde{\nabla}f(x_s)^2} \right) \\ 4: & x_{t+1} \leftarrow \operatorname*{arg\,min}_{x \in \mathcal{C}} \eta \langle x, \tilde{\nabla}f(x_t) \rangle + \frac{1}{2} \|x - x_t\|_{H_t}^2 \end{array}$ 

Simultaneously proposed by Duchi et al. (2011) and McMahan & Streeter (2010):

AlgorithmAdaptive Gradient (AdaGrad)Input: $x_0 \in C, \ \delta > 0, \ \eta > 0$ 1:for t = 0 to T - 1 do2:Update the gradient estimator  $\tilde{\nabla}f(x_t)$ 3: $H_t \leftarrow \text{diag}\left(\delta 1 + \sqrt{\sum_{s=0}^t \tilde{\nabla}f(x_s)^2}\right)$ 4: $x_{t+1} \leftarrow \underset{x \in C}{\arg \min \eta \langle x, \tilde{\nabla}f(x_t) \rangle} + \frac{1}{2} ||x - x_t||_{H_t}^2$ 

Simultaneously proposed by Duchi et al. (2011) and McMahan & Streeter (2010):

AlgorithmAdaptive Gradient (AdaGrad)Input: $x_0 \in C, \ \delta > 0, \ \eta > 0$ 1:for t = 0 to T - 1 do2:Update the gradient estimator  $\tilde{\nabla}f(x_t)$ 3: $H_t \leftarrow \text{diag}\left(\delta 1 + \sqrt{\sum_{s=0}^t \tilde{\nabla}f(x_s)^2}\right)$ 4: $x_{t+1} \leftarrow \underset{x \in C}{\arg\min \eta \langle x, \tilde{\nabla}f(x_t) \rangle} + \frac{1}{2} ||x - x_t||_{H_t}^2$ 

• We denote  $||u||_{H_t}^2 = \langle u, H_t u \rangle$ 

Simultaneously proposed by Duchi et al. (2011) and McMahan & Streeter (2010):

AlgorithmAdaptive Gradient (AdaGrad)Input: $x_0 \in C, \ \delta > 0, \ \eta > 0$ 1:for t = 0 to T - 1 do2:Update the gradient estimator  $\tilde{\nabla}f(x_t)$ 3: $H_t \leftarrow \text{diag}\left(\delta 1 + \sqrt{\sum_{s=0}^t \tilde{\nabla}f(x_s)^2}\right)$ 4: $x_{t+1} \leftarrow \arg\min_{x \in C} \eta \langle x, \tilde{\nabla}f(x_t) \rangle + \frac{1}{2} ||x - x_t||^2_{H_t}$ 

• We denote  $||u||_{H_t}^2 = \langle u, H_t u \rangle$ 

• The default value for the offset is  $\delta \leftarrow 10^{-8}$ 

Simultaneously proposed by Duchi et al. (2011) and McMahan & Streeter (2010):

AlgorithmAdaptive Gradient (AdaGrad)Input:  $x_0 \in C, \ \delta > 0, \ \eta > 0$ 1: for t = 0 to T - 1 do2: Update the gradient estimator  $\tilde{\nabla} f(x_t)$ 3:  $H_t \leftarrow \text{diag}\left(\delta 1 + \sqrt{\sum_{s=0}^t \tilde{\nabla} f(x_s)^2}\right)$ 4:  $x_{t+1} \leftarrow \underset{x \in C}{\arg \min \eta \langle x, \tilde{\nabla} f(x_t) \rangle} + \frac{1}{2} ||x - x_t||_{H_t}^2$ 

• We denote  $||u||_{H_t}^2 = \langle u, H_t u \rangle$ 

• The default value for the offset is  $\delta \leftarrow 10^{-8}$ 

We can rewrite

$$x_{t+1} \leftarrow \operatorname*{arg\,min}_{x \in \mathcal{C}} \|x - (x_t - \eta H_t^{-1} \tilde{\nabla} f(x_t))\|_{H_t}$$

We can rewrite

$$x_{t+1} \leftarrow \operatorname*{arg\,min}_{x \in \mathcal{C}} \|x - (x_t - \eta H_t^{-1} \tilde{\nabla} f(x_t))\|_{H_t}$$

Ignoring the constraint set  $\ensuremath{\mathcal{C}}$  for ease of exposition, we obtain

$$x_{t+1} \leftarrow x_t - \eta H_t^{-1} \tilde{\nabla} f(x_t)$$

We can rewrite

$$x_{t+1} \leftarrow \operatorname*{arg\,min}_{x \in \mathcal{C}} \|x - (x_t - \eta H_t^{-1} \tilde{\nabla} f(x_t))\|_{H_t}$$

Ignoring the constraint set  $\ensuremath{\mathcal{C}}$  for ease of exposition, we obtain

$$x_{t+1} \leftarrow x_t - \eta H_t^{-1} \tilde{\nabla} f(x_t)$$

i.e., for every feature  $i \in \llbracket 1, n \rrbracket$ ,

$$[\mathbf{x}_{t+1}]_i \leftarrow [\mathbf{x}_t]_i - \frac{\eta[\tilde{\nabla}f(\mathbf{x}_t)]_i}{\delta + \sqrt{\sum_{s=0}^t [\tilde{\nabla}f(\mathbf{x}_s)]_i^2}}$$

We can rewrite

$$x_{t+1} \leftarrow \operatorname*{arg\,min}_{x \in \mathcal{C}} \|x - (x_t - \eta H_t^{-1} \tilde{\nabla} f(x_t))\|_{H_t}$$

Ignoring the constraint set  $\ensuremath{\mathcal{C}}$  for ease of exposition, we obtain

$$x_{t+1} \leftarrow x_t - \eta H_t^{-1} \tilde{\nabla} f(x_t)$$

i.e., for every feature  $i \in \llbracket 1, n \rrbracket$ ,

$$[\mathbf{x}_{t+1}]_i \leftarrow [\mathbf{x}_t]_i - \frac{\eta[\tilde{\nabla}f(\mathbf{x}_t)]_i}{\delta + \sqrt{\sum_{s=0}^t [\tilde{\nabla}f(\mathbf{x}_s)]_i^2}}$$

• The offset  $\delta$  prevents from dividing by zero

We can rewrite

$$x_{t+1} \leftarrow \operatorname*{arg\,min}_{x \in \mathcal{C}} \|x - (x_t - \eta H_t^{-1} \tilde{\nabla} f(x_t))\|_{H_t}$$

Ignoring the constraint set  $\ensuremath{\mathcal{C}}$  for ease of exposition, we obtain

$$x_{t+1} \leftarrow x_t - \eta H_t^{-1} \tilde{\nabla} f(x_t)$$

i.e., for every feature  $i \in \llbracket 1, n \rrbracket$ ,

$$[\mathbf{x}_{t+1}]_i \leftarrow [\mathbf{x}_t]_i - \frac{\eta[\tilde{\nabla}f(\mathbf{x}_t)]_i}{\delta + \sqrt{\sum_{s=0}^t [\tilde{\nabla}f(\mathbf{x}_s)]_i^2}}$$

- The offset  $\delta$  prevents from dividing by zero
- The step-size automatically adjusts to the geometry of the problem

We can rewrite

$$x_{t+1} \leftarrow \operatorname*{arg\,min}_{x \in \mathcal{C}} \|x - (x_t - \eta H_t^{-1} \tilde{\nabla} f(x_t))\|_{H_t}$$

Ignoring the constraint set  $\ensuremath{\mathcal{C}}$  for ease of exposition, we obtain

$$x_{t+1} \leftarrow x_t - \eta H_t^{-1} \tilde{\nabla} f(x_t)$$

i.e., for every feature  $i \in \llbracket 1, n \rrbracket$ ,

$$[\mathbf{x}_{t+1}]_i \leftarrow [\mathbf{x}_t]_i - \frac{\eta[\tilde{\nabla}f(\mathbf{x}_t)]_i}{\delta + \sqrt{\sum_{s=0}^t [\tilde{\nabla}f(\mathbf{x}_s)]_i^2}}$$

- The offset  $\delta$  prevents from dividing by zero
- The step-size automatically adjusts to the geometry of the problem
- Larger step-sizes are given to infrequent (but potentially very informative) features whenever they appear so that they do not go unnoticed


















• How can we use adaptive gradients in FW?

- How can we use adaptive gradients in FW?
- Let  $G_t = H_t^{-1} \tilde{\nabla} f(x_t)$ , then unconstrained AdaGrad is

$$x_{t+1} \leftarrow x_t - \eta G_t$$

- How can we use adaptive gradients in FW?
- Let  $G_t = H_t^{-1} \tilde{\nabla} f(x_t)$ , then unconstrained AdaGrad is

$$x_{t+1} \leftarrow x_t - \eta G_t$$

Could we do

$$v_t \leftarrow \operatorname*{arg\,min}_{v \in \mathcal{C}} \langle v, G_t \rangle$$
$$x_{t+1} \leftarrow x_t + \gamma_t (v_t - x_t)$$

as did FW for unconstrained gradient descent (for which  $G_t = \nabla f(x_t)$ )?

- How can we use adaptive gradients in FW?
- Let  $G_t = H_t^{-1} \tilde{\nabla} f(x_t)$ , then unconstrained AdaGrad is

$$x_{t+1} \leftarrow x_t - \eta G_t$$

Could we do

$$v_t \leftarrow \operatorname*{arg\,min}_{v \in \mathcal{C}} \langle v, G_t \rangle$$
$$x_{t+1} \leftarrow x_t + \gamma_t (v_t - x_t)$$

as did FW for unconstrained gradient descent (for which  $G_t = \nabla f(x_t)$ )?

• We would likely lose the precious properties of the descent directions of AdaGrad

• Instead, consider the constrained subproblem occurring at every iteration:

$$x_{t+1} \leftarrow \operatorname*{arg\,min}_{x \in \mathcal{C}} \eta \langle x, \tilde{
abla} f(x_t) 
angle + rac{1}{2} \|x - x_t\|_{H_t}^2$$

• Instead, consider the constrained subproblem occurring at every iteration:

$$x_{t+1} \leftarrow rgmin_{x \in \mathcal{C}} \eta \langle x, ilde{
abla} f(x_t) 
angle + rac{1}{2} \|x - x_t\|_{H_t}^2$$

• This can become quite expensive and inefficient overall (AdaGrad is usually used for unconstrained optimization)

• Instead, consider the constrained subproblem occurring at every iteration:

$$x_{t+1} \leftarrow rgmin_{x \in \mathcal{C}} \eta \langle x, ilde{
abla} f(x_t) 
angle + rac{1}{2} \|x - x_t\|_{H_t}^2$$

• This can become quite expensive and inefficient overall (AdaGrad is usually used for unconstrained optimization)

Idea:

• Solve the subproblem using FW (sliding technique of Lan & Zhou (2016))

• Instead, consider the constrained subproblem occurring at every iteration:

$$egin{aligned} x_{t+1} \leftarrow rgmin_{x\in\mathcal{C}} \eta\langle x, ilde{
abla}f(x_t)
angle + rac{1}{2}\|x-x_t\|_{\mathcal{H}_t}^2 \end{aligned}$$

• This can become quite expensive and inefficient overall (AdaGrad is usually used for unconstrained optimization)

Idea:

- Solve the subproblem using FW (sliding technique of Lan & Zhou (2016))
- Run only a small and fixed number K of iterations of FW ( $K \sim 5$ )

• Instead, consider the constrained subproblem occurring at every iteration:

$$egin{aligned} x_{t+1} \leftarrow rgmin_{x\in\mathcal{C}} \eta\langle x, ilde{
abla}f(x_t)
angle + rac{1}{2}\|x-x_t\|_{\mathcal{H}_t}^2 \end{aligned}$$

• This can become quite expensive and inefficient overall (AdaGrad is usually used for unconstrained optimization)

Idea:

- Solve the subproblem using FW (sliding technique of Lan & Zhou (2016))
- Run only a small and fixed number K of iterations of FW ( $K \sim 5$ )
- That is, we claim that leveraging just a small amount of information from the adaptive metric  $H_t$  is enough

#### **Template** Frank-Wolfe with adaptive gradients

**Input:**  $x_0 \in \mathcal{C}, \ 0 < \lambda_t^- \leq \lambda_{t+1}^- \leq \lambda_{t+1}^+ \leq \lambda_t^+, \ K \in \mathbb{N} \setminus \{0\}, \ \eta > 0, \ \gamma_t \in [0, 1]$ 1: for t = 0 to T - 1 do Update the gradient estimator  $\tilde{\nabla} f(x_t)$ 2: 3: Update the diagonal matrix  $H_t$  and clip its entries to  $[\lambda_t^-, \lambda_t^+]$  $v_0^{(t)} \leftarrow x_t$ 4: 5: for k = 0 to K - 1 do  $\nabla Q_t(y_k^{(t)}) \leftarrow \tilde{\nabla} f(x_t) + \frac{1}{n_t} H_t(y_k^{(t)} - x_t)$ 6:  $v_{\iota}^{(t)} \leftarrow \arg\min\langle \nabla Q_t(y_k^{(t)}), v \rangle$ 7:  $\gamma_k^{(t)} \leftarrow \min\left\{\eta_t \frac{\langle \nabla Q_t(y_k^{(t)}), y_k^{(t)} - v_k^{(t)} \rangle}{\|v_k^{(t)} - v_t^{(t)}\|_{L^1}^2}, \gamma_t\right\}$ 8:  $y_{k+1} \leftarrow y_k^{(t)} + \gamma_k^{(t)} (v_k^{(t)} - y_k^{(t)})$  $x_{t+1} \leftarrow y_k^{(t)}$ 9: 10:

#### **Template** Frank-Wolfe with adaptive gradients

**Input:**  $x_0 \in \mathcal{C}, \ 0 < \lambda_t^- \leq \lambda_{t+1}^- \leq \lambda_{t+1}^+ \leq \lambda_t^+, \ K \in \mathbb{N} \setminus \{0\}, \ \eta > 0, \ \gamma_t \in [0, 1]$ 1: for t = 0 to T - 1 do Update the gradient estimator  $\tilde{\nabla} f(x_t)$ 2: 3: Update the diagonal matrix  $H_t$  and clip its entries to  $[\lambda_t^-, \lambda_t^+]$  $v_0^{(t)} \leftarrow x_t$ 4: 5: for k = 0 to K - 1 do  $\nabla Q_t(y_k^{(t)}) \leftarrow \tilde{\nabla} f(x_t) + \frac{1}{n_t} H_t(y_k^{(t)} - x_t)$ 6:  $v_{\iota}^{(t)} \leftarrow \arg\min\langle \nabla Q_t(y_k^{(t)}), v \rangle$ 7:  $\gamma_k^{(t)} \leftarrow \min\left\{\eta_t \frac{\langle \nabla Q_t(y_k^{(t)}), y_k^{(t)} - v_k^{(t)} \rangle}{\|v_k^{(t)} - v_t^{(t)}\|_{L^1}^2}, \gamma_t\right\}$ 8:

9:  $y_{k+1} \leftarrow y_k^{(t)} + \gamma_k^{(t)} (v_k^{(t)} - y_k^{(t)})$ 10:  $x_{t+1} \leftarrow y_k^{(t)}$ 

#### **Template** Frank-Wolfe with adaptive gradients

**Input:**  $x_0 \in \mathcal{C}, \ 0 < \lambda_t^- \leq \lambda_{t+1}^- \leq \lambda_{t+1}^+ \leq \lambda_t^+, \ K \in \mathbb{N} \setminus \{0\}, \ \eta > 0, \ \gamma_t \in [0, 1]$ 1: for t = 0 to T - 1 do Update the gradient estimator  $\tilde{\nabla} f(x_t)$ 2: 3: Update the diagonal matrix  $H_t$  and clip its entries to  $[\lambda_t^-, \lambda_t^+]$  $v_0^{(t)} \leftarrow x_t$ 4: 5: for k = 0 to K - 1 do  $\nabla Q_t(y_k^{(t)}) \leftarrow \tilde{\nabla} f(x_t) + \frac{1}{n_t} H_t(y_k^{(t)} - x_t)$ 6:  $v_{\iota}^{(t)} \leftarrow \arg\min\langle \nabla Q_t(y_k^{(t)}), v \rangle$ 7:  $\gamma_k^{(t)} \leftarrow \min\left\{\eta_t \frac{\langle \nabla Q_t(y_k^{(t)}), y_k^{(t)} - v_k^{(t)} \rangle}{\|v_k^{(t)} - v_t^{(t)}\|_{L^1}^2}, \gamma_t\right\}$ 8:  $y_{k+1} \leftarrow y_k^{(t)} + \gamma_k^{(t)} (v_k^{(t)} - y_k^{(t)})$  $x_{t+1} \leftarrow y_k^{(t)}$ 9: 10:

### **Template** Frank-Wolfe with adaptive gradients

Input:  $x_0 \in C$ ,  $0 < \lambda_t^- \leq \lambda_{t+1}^- \leq \lambda_{t+1}^+ \leq \lambda_t^+$ ,  $K \in \mathbb{N} \setminus \{0\}$ ,  $\eta > 0$ ,  $\gamma_t \in [0, 1]$ 1: for t = 0 to T - 1 do 2: Update the gradient estimator  $\tilde{\nabla}f(x_t)$ 

- 3: Update the diagonal matrix  $H_t$  and clip its entries to  $[\lambda_t^-, \lambda_t^+]$
- 4:  $y_0^{(t)} \leftarrow x_t$
- 5: for k = 0 to K 1 do

6: 
$$\nabla Q_t(y_k^{(t)}) \leftarrow \tilde{\nabla} f(x_t) + \frac{1}{\eta_t} H_t(y_k^{(t)} - x_t)$$

7: 
$$v_k^{(t)} \leftarrow \underset{v \in \mathcal{C}}{\operatorname{arg\,min}} \langle \nabla Q_t(y_k^{(t)}), v \rangle$$

8: 
$$\gamma_k^{(t)} \leftarrow \min\left\{\eta_t \frac{\langle \nabla Q_t(y_k^{(t)}), y_k^{(t)} - v_k^{(t)} \rangle}{\|y_k^{(t)} - v_k^{(t)}\|_{H_t}^2}, \gamma_t\right\}$$

9: 
$$y_{k+1} \leftarrow y_k^{(t)} + \gamma_k^{(t)} (v_k^{(t)} - y_k^{(t)})$$

10:  $x_{t+1} \leftarrow y_K^{(t)}$ 

• Lines 4–9 apply 
$$K$$
 iterations of FW to  

$$\min_{x \in \mathcal{C}} \left\{ Q_t(x) \coloneqq f(x_t) + \langle x - x_t, \tilde{\nabla} f(x_t) \rangle + \frac{1}{2\eta_t} \| x - x_t \|_{H_t}^2 \right\}$$

#### **Template** Frank-Wolfe with adaptive gradients

**Input:**  $x_0 \in \mathcal{C}, \ 0 < \lambda_t^- \leq \lambda_{t+1}^- \leq \lambda_{t+1}^+ \leq \lambda_t^+, \ K \in \mathbb{N} \setminus \{0\}, \ \eta > 0, \ \gamma_t \in [0, 1]$ 1: for t = 0 to T - 1 do Update the gradient estimator  $\tilde{\nabla} f(x_t)$ 2: 3: Update the diagonal matrix  $H_t$  and clip its entries to  $[\lambda_t^-, \lambda_t^+]$  $v_0^{(t)} \leftarrow x_t$ 4: 5: for k = 0 to K - 1 do  $abla Q_t(y_k^{(t)}) \leftarrow ilde{
abla} f(x_t) + rac{1}{n_t} H_t(y_k^{(t)} - x_t)$ 6:  $v_k^{(t)} \leftarrow \arg\min\langle \nabla Q_t(y_k^{(t)}), v \rangle$ 7:  $\gamma_k^{(t)} \leftarrow \min\left\{\eta_t \frac{\langle \nabla Q_t(y_k^{(t)}), y_k^{(t)} - v_k^{(t)} \rangle}{\|v_k^{(t)} - v_t^{(t)}\|_{i}^2}, \gamma_t\right\}$ 8:  $y_{k+1} \leftarrow y_k^{(t)} + \gamma_k^{(t)} (v_k^{(t)} - y_k^{(t)})$  $x_{t+1} \leftarrow y_{\nu}^{(t)}$ 9. 10:

• Lines 4–9 apply K iterations of FW to  $\min_{x \in \mathcal{C}} \left\{ Q_t(x) \coloneqq f(x_t) + \langle x - x_t, \tilde{\nabla} f(x_t) \rangle + \frac{1}{2\eta_t} \| x - x_t \|_{H_t}^2 \right\}$ 

### **Template** Frank-Wolfe with adaptive gradients

**Input:**  $x_0 \in \mathcal{C}, \ 0 < \lambda_t^- \leq \lambda_{t+1}^- \leq \lambda_{t+1}^+ \leq \lambda_t^+, \ K \in \mathbb{N} \setminus \{0\}, \ \eta > 0, \ \gamma_t \in [0, 1]$ 1: for t = 0 to T - 1 do Update the gradient estimator  $\tilde{\nabla} f(x_t)$ 2: 3: Update the diagonal matrix  $H_t$  and clip its entries to  $[\lambda_t^-, \lambda_t^+]$  $v_0^{(t)} \leftarrow x_t$ 4: 5: for k = 0 to K - 1 do  $\nabla Q_t(y_k^{(t)}) \leftarrow \tilde{\nabla} f(x_t) + \frac{1}{n_t} H_t(y_k^{(t)} - x_t)$ 6:  $v_k^{(t)} \leftarrow \arg\min\langle \nabla Q_t(y_k^{(t)}), v \rangle$ 7:  $\gamma_k^{(t)} \leftarrow \min\left\{\eta_t \frac{\langle \nabla Q_t(y_k^{(t)}), y_k^{(t)} - v_k^{(t)} \rangle}{\|v_t^{(t)} - v_t^{(t)}\|_{\cdot}^2}, \gamma_t\right\}$ 8:  $y_{k+1} \leftarrow y_k^{(t)} + \gamma_k^{(t)} (v_k^{(t)} - y_k^{(t)})$  $x_{t+1} \leftarrow y_{\kappa}^{(t)}$ 9. 10:

- Lines 4–9 apply K iterations of FW to  $\min_{x \in \mathcal{C}} \left\{ Q_t(x) \coloneqq f(x_t) + \langle x - x_t, \tilde{\nabla} f(x_t) \rangle + \frac{1}{2\eta_t} \| x - x_t \|_{H_t}^2 \right\}$
- Define AdaX depending on the strategy for  $\tilde{\nabla} f(x_t)$ : AdaSFW, AdaSVRF, etc.

#### Theorem

Let  $f_1, \ldots, f_m \colon \mathbb{R}^n \to \mathbb{R}$  be L-smooth convex functions with  $G = \max_{\mathcal{C}} \|\nabla f_i\|_2$ . Then AdaSFW with  $b_t \leftarrow (G(t+2)/(LD))^2$ ,  $\eta_t \leftarrow \lambda_t^-/L$ , and  $\gamma_t \leftarrow 2/(t+2)$  satisfies

$$\mathbb{E}[f(x_t)] - \min_{\mathcal{C}} f \leqslant \frac{2LD^2(K+1+\kappa)}{t+1}$$

where  $\kappa = \lambda_0^+ / \lambda_0^-$ 

#### Theorem

Let  $f_1, \ldots, f_m \colon \mathbb{R}^n \to \mathbb{R}$  be L-smooth convex functions with  $G = \max_{\mathcal{C}} \|\nabla f_i\|_2$ . Then AdaSFW with  $b_t \leftarrow (G(t+2)/(LD))^2$ ,  $\eta_t \leftarrow \lambda_t^-/L$ , and  $\gamma_t \leftarrow 2/(t+2)$  satisfies

$$\mathbb{E}[f(x_t)] - \min_{\mathcal{C}} f \leqslant \frac{2LD^2(K+1+\kappa)}{t+1}$$

where  $\kappa = \lambda_0^+ / \lambda_0^-$ 

• In practice, no need to know G,L,D and simply set  $b_t=\Theta(t^2)$ 

#### Theorem

Let  $f_1, \ldots, f_m \colon \mathbb{R}^n \to \mathbb{R}$  be L-smooth convex functions with  $G = \max_{\mathcal{C}} \|\nabla f_i\|_2$ . Then AdaSFW with  $b_t \leftarrow (G(t+2)/(LD))^2$ ,  $\eta_t \leftarrow \lambda_t^-/L$ , and  $\gamma_t \leftarrow 2/(t+2)$  satisfies

$$\mathbb{E}[f(x_t)] - \min_{\mathcal{C}} f \leqslant \frac{2LD^2(K+1+\kappa)}{t+1}$$

where  $\kappa = \lambda_0^+ / \lambda_0^-$ 

- In practice, no need to know G, L, D and simply set  $b_t = \Theta(t^2)$
- AdaSVRF and AdaCSFW also converge at a rate  $\mathcal{O}(1/t)$  with  $b_t \sim t$  and  $b_t \sim 1$  respectively

#### Theorem

Let  $f_1, \ldots, f_m \colon \mathbb{R}^n \to \mathbb{R}$  be L-smooth convex functions with  $G = \max_{\mathcal{C}} \|\nabla f_i\|_2$ . Then AdaSFW with  $b_t \leftarrow (G(t+2)/(LD))^2$ ,  $\eta_t \leftarrow \lambda_t^-/L$ , and  $\gamma_t \leftarrow 2/(t+2)$  satisfies

$$\mathbb{E}[f(x_t)] - \min_{\mathcal{C}} f \leqslant \frac{2LD^2(K+1+\kappa)}{t+1}$$

where  $\kappa = \lambda_0^+ / \lambda_0^-$ 

- In practice, no need to know G, L, D and simply set  $b_t = \Theta(t^2)$
- AdaSVRF and AdaCSFW also converge at a rate  $\mathcal{O}(1/t)$  with  $b_t \sim t$  and  $b_t \sim 1$  respectively
- If  $f_1, \ldots, f_m$  are nonconvex, then AdaSFW converges to a stationary point at a rate  $\mathcal{O}(1/t^{1/2-\nu})$  with  $b_t \sim t$  and  $\gamma_t \leftarrow 1/(t+1)^{1/2+\nu}$ ,  $0 < \nu < 1/2$

• We compare our method to SFW, SVRF, SPIDER-FW, ORGFW, and CSFW, as well as AdaGrad and AMSGrad on a wide range of experiments

- We compare our method to SFW, SVRF, SPIDER-FW, ORGFW, and CSFW, as well as AdaGrad and AMSGrad on a wide range of experiments
- In addition, we run AdamSFW, a variant of AdaSFW with momentum inspired by Kingma & Ba (2015); Reddi et al. (2018)

- We compare our method to SFW, SVRF, SPIDER-FW, ORGFW, and CSFW, as well as AdaGrad and AMSGrad on a wide range of experiments
- In addition, we run AdamSFW, a variant of AdaSFW with momentum inspired by Kingma & Ba (2015); Reddi et al. (2018)
- We set  $K \sim 5$

# Convex objectives

• Support vector classification on a synthetic dataset





 Linear regression on the YearPredictionMSD dataset





Logistic regression on the RCV1 dataset



# Nonconvex objectives

10<sup>1</sup>

sol principal to a sol principal to a sol principal to a solution of the solut

88

84

82

80

fest accuracy (%)

- Neural network with one hidden • layer on the IMDB dataset
- AdaGrad AdaGrad AMSGrad AMSGrad 1.5 SFW SEW  $10^{\circ}$ Training loss SVRF SVRF SPIDER-FW SPIDER-FW ORGEW ORGFW AdaSFW AdaSFW AdamSFW AdamSFW  $10^{-2}$ 0.3 82 80 76 76 76 70 80 76 76 80 76 72 70 80 80 76 86 Test 66 64 60 ò 25 50 75 100 10 20 0 5 10 15 20 0 100 200 0 Epoch CPU time (h) Epoch CPU time (s)

•

Convolutional network on the

CIFAR-10 dataset

Each layer is constrained into an  $\ell_{\infty}$ -ball •

Let  $C \subset \mathbb{R}^n$  be a compact convex set and  $x^* \in C$ . Then  $x^*$  can be represented as a convex combination of at most n + 1 vertices of C.

 In R<sup>2</sup>, every point in C is a convex combination of at most 3 vertices



Let  $C \subset \mathbb{R}^n$  be a compact convex set and  $x^* \in C$ . Then  $x^*$  can be represented as a convex combination of at most n + 1 vertices of C.

In R<sup>2</sup>, every point in C is a convex combination of at most 3 vertices



Let  $C \subset \mathbb{R}^n$  be a compact convex set and  $x^* \in C$ . Then  $x^*$  can be represented as a convex combination of at most n + 1 vertices of C.

In R<sup>2</sup>, every point in C is a convex combination of at most 3 vertices



- In R<sup>2</sup>, every point in C is a convex combination of at most 3 vertices
- Can we reduce *n* + 1 when we can afford an *ε*-approximation?



- In R<sup>2</sup>, every point in C is a convex combination of at most 3 vertices
- Can we reduce *n* + 1 when we can afford an *ε*-approximation?



- In R<sup>2</sup>, every point in C is a convex combination of at most 3 vertices
- Can we reduce *n* + 1 when we can afford an *ε*-approximation?


### Theorem (Carathéodory, 1907)

Let  $C \subset \mathbb{R}^n$  be a compact convex set and  $x^* \in C$ . Then  $x^*$  can be represented as a convex combination of at most n + 1 vertices of C.

- In R<sup>2</sup>, every point in C is a convex combination of at most 3 vertices
- Can we reduce *n* + 1 when we can afford an *ε*-approximation?



### Theorem (Carathéodory, 1907)

Let  $C \subset \mathbb{R}^n$  be a compact convex set and  $x^* \in C$ . Then  $x^*$  can be represented as a convex combination of at most n + 1 vertices of C.

- In R<sup>2</sup>, every point in C is a convex combination of at most 3 vertices
- Can we reduce *n* + 1 when we can afford an *ε*-approximation?
- Define the *cardinality* of x ∈ C as the number of vertices in a given convex decomposition of x



### Problem

Find  $x \in C$  with low cardinality satisfying  $||x - x^*||_p \leq \varepsilon$ .

### Problem

Find  $x \in C$  with low cardinality satisfying  $||x - x^*||_p \leq \varepsilon$ .

• Applications in game theory, combinatorial optimization, etc.

### Problem

Find  $x \in C$  with low cardinality satisfying  $||x - x^*||_p \leq \varepsilon$ .

• Applications in game theory, combinatorial optimization, etc.

### Theorem (Barman, 2015)

### Problem

Find  $x \in C$  with low cardinality satisfying  $||x - x^*||_p \leq \varepsilon$ .

• Applications in game theory, combinatorial optimization, etc.

### Theorem (Barman, 2015)

Let  $p \in [2, +\infty[$ . Then there exists  $x \in C$  with cardinality  $\mathcal{O}(pD_p^2/\varepsilon^2)$  satisfying  $||x - x^*||_p \leq \varepsilon$ , where  $D_p$  is the diameter of C in the  $\ell_p$ -norm.

• This result is independent of the ambient dimension *n* 

### Problem

Find  $x \in C$  with low cardinality satisfying  $||x - x^*||_p \leq \varepsilon$ .

• Applications in game theory, combinatorial optimization, etc.

### Theorem (Barman, 2015)

- This result is independent of the ambient dimension *n*
- The bound is tight (Mirrokni et al., 2017)

### Problem

Find  $x \in C$  with low cardinality satisfying  $||x - x^*||_p \leq \varepsilon$ .

• Applications in game theory, combinatorial optimization, etc.

### Theorem (Barman, 2015)

- This result is independent of the ambient dimension *n*
- The bound is tight (Mirrokni et al., 2017)
- Probabilistic proofs by Pisier (1981); Barman (2015)

### Problem

Find  $x \in C$  with low cardinality satisfying  $||x - x^*||_p \leq \varepsilon$ .

• Applications in game theory, combinatorial optimization, etc.

### Theorem (Barman, 2015)

- This result is independent of the ambient dimension *n*
- The bound is tight (Mirrokni et al., 2017)
- Probabilistic proofs by Pisier (1981); Barman (2015)
- Deterministic proof by Mirrokni et al. (2017) using mirror descent (Nemirovsky & Yudin, 1983) on the dual problem

### Problem

Find  $x \in C$  with low cardinality satisfying  $||x - x^*||_p \leq \varepsilon$ .

• Applications in game theory, combinatorial optimization, etc.

### Theorem (Barman, 2015)

- This result is independent of the ambient dimension *n*
- The bound is tight (Mirrokni et al., 2017)
- Probabilistic proofs by Pisier (1981); Barman (2015)
- Deterministic proof by Mirrokni et al. (2017) using mirror descent (Nemirovsky & Yudin, 1983) on the dual problem
- Can we solve  $\min_{x \in C} ||x x^*||_p$  by sequentially picking up vertices?

Let  $p \in [2, +\infty[$  and  $f(x) = \frac{1}{2} ||x - x^*||_p^2$ . Then f is convex, (p-1)-smooth, and 1-gradient dominated w.r.t. the  $\ell_p$ -norm.

Let  $p \in [2, +\infty[$  and  $f(x) = \frac{1}{2} ||x - x^*||_p^2$ . Then f is convex, (p-1)-smooth, and 1-gradient dominated w.r.t. the  $\ell_p$ -norm.

• If  $p \in [2, +\infty[$ , run FW on  $\min_{x \in C} \frac{1}{2} ||x - x^*||_p^2$  and count the number of iterations to reach  $\varepsilon^2/2$ -convergence

Let  $p \in [2, +\infty[$  and  $f(x) = \frac{1}{2} ||x - x^*||_p^2$ . Then f is convex, (p-1)-smooth, and 1-gradient dominated w.r.t. the  $\ell_p$ -norm.

• If  $p \in [2, +\infty[$ , run FW on  $\min_{x \in C} \frac{1}{2} ||x - x^*||_p^2$  and count the number of iterations to reach  $\varepsilon^2/2$ -convergence

#### Lemma

Let  $p \in [1, 2[ \cup \{+\infty\} \text{ and } f(x) = ||x - x^*||_p$ . Then f is convex and Lipschitz continuous w.r.t. the  $\ell_2$ -norm, with constant  $n^{1/p-1/2}$  if  $p \in [1, 2[$ , else 1 if  $p = +\infty$ .

Let  $p \in [2, +\infty[$  and  $f(x) = \frac{1}{2} ||x - x^*||_p^2$ . Then f is convex, (p-1)-smooth, and 1-gradient dominated w.r.t. the  $\ell_p$ -norm.

• If  $p \in [2, +\infty[$ , run FW on  $\min_{x \in C} \frac{1}{2} ||x - x^*||_p^2$  and count the number of iterations to reach  $\varepsilon^2/2$ -convergence

#### Lemma

Let  $p \in [1, 2[ \cup \{+\infty\} \text{ and } f(x) = ||x - x^*||_p$ . Then f is convex and Lipschitz continuous w.r.t. the  $\ell_2$ -norm, with constant  $n^{1/p-1/2}$  if  $p \in [1, 2[$ , else 1 if  $p = +\infty$ .

• If  $p \in [1, 2[ \cup \{+\infty\}, \text{ run HCGS on } \min_{x \in C} ||x - x^*||_p$  and count the number of iterations to reach  $\varepsilon$ -convergence

Smoothen the problem (Argyriou et al., 2014) by taking the Moreau envelope (Moreau, 1965):

Smoothen the problem (Argyriou et al., 2014) by taking the Moreau envelope (Moreau, 1965):

$$f_{\beta}(x) = \min_{y \in \mathbb{R}^n} f(y) + \frac{1}{2\beta} ||x - y||_2^2$$

Smoothen the problem (Argyriou et al., 2014) by taking the Moreau envelope (Moreau, 1965):

$$f_{\beta}(x) = \min_{y \in \mathbb{R}^n} f(y) + \frac{1}{2\beta} \|x - y\|_2^2 \quad \text{and} \quad \nabla f_{\beta}(x) = \frac{1}{\beta} (x - \operatorname{prox}_{\beta f}(x))$$

Smoothen the problem (Argyriou et al., 2014) by taking the Moreau envelope (Moreau, 1965):

$$f_{eta}(x) = \min_{y \in \mathbb{R}^n} f(y) + rac{1}{2eta} \|x - y\|_2^2 \quad ext{and} \quad 
abla f_{eta}(x) = rac{1}{eta} (x - ext{prox}_{eta f}(x))$$

Algorithm Hybrid Conditional Gradient-Smoothing (HCGS)

Input: 
$$x_0 \in C$$
,  $G_2, D_2$   
1: for  $t = 0$  to  $T - 1$  do  
2:  $\beta_t \leftarrow 2(D_2/G_2)/\sqrt{t+2}$   
3:  $v_t \leftarrow \arg \min_{v \in C} \langle v, \nabla f_{\beta_t}(x_t) \rangle$   
4:  $\gamma_t \leftarrow 2/(t+2)$   
5:  $x_{t+1} \leftarrow x_t + \gamma_t(v_t - x_t)$ 

Smoothen the problem (Argyriou et al., 2014) by taking the Moreau envelope (Moreau, 1965):

$$f_{eta}(x) = \min_{y \in \mathbb{R}^n} f(y) + rac{1}{2eta} \|x - y\|_2^2 \quad ext{and} \quad 
abla f_{eta}(x) = rac{1}{eta} (x - ext{prox}_{eta f}(x))$$

Algorithm Hybrid Conditional Gradient-Smoothing (HCGS)

Input:  $x_0 \in C$ ,  $G_2$ ,  $D_2$ 1: for t = 0 to T - 1 do 2:  $\beta_t \leftarrow 2(D_2/G_2)/\sqrt{t+2}$ 3:  $v_t \leftarrow \arg\min_{v \in C} \langle v, \nabla f_{\beta_t}(x_t) \rangle$ 4:  $\gamma_t \leftarrow 2/(t+2)$ 5:  $x_{t+1} \leftarrow x_t + \gamma_t(v_t - x_t)$ 

Smoothen the problem (Argyriou et al., 2014) by taking the Moreau envelope (Moreau, 1965):

$$f_{eta}(x) = \min_{y \in \mathbb{R}^n} f(y) + rac{1}{2eta} \|x - y\|_2^2 \quad ext{and} \quad 
abla f_{eta}(x) = rac{1}{eta} (x - ext{prox}_{eta f}(x))$$

Algorithm Hybrid Conditional Gradient-Smoothing (HCGS)

Input:  $x_0 \in C$ ,  $G_2$ ,  $D_2$ 1: for t = 0 to T - 1 do 2:  $\beta_t \leftarrow 2(D_2/G_2)/\sqrt{t+2}$ 3:  $v_t \leftarrow \arg\min_{v \in C} \langle v, \nabla f_{\beta_t}(x_t) \rangle$ 4:  $\gamma_t \leftarrow 2/(t+2)$ 5:  $x_{t+1} \leftarrow x_t + \gamma_t(v_t - x_t)$ 

Smoothen the problem (Argyriou et al., 2014) by taking the Moreau envelope (Moreau, 1965):

$$f_{eta}(x) = \min_{y \in \mathbb{R}^n} f(y) + rac{1}{2eta} \|x - y\|_2^2 \quad ext{and} \quad 
abla f_{eta}(x) = rac{1}{eta} (x - ext{prox}_{eta f}(x))$$

Algorithm Hybrid Conditional Gradient-Smoothing (HCGS)

- Input:  $x_0 \in C$ ,  $G_2$ ,  $D_2$ 1: for t = 0 to T - 1 do 2:  $\beta_t \leftarrow 2(D_2/G_2)/\sqrt{t+2}$ 3:  $v_t \leftarrow \arg\min_{v \in C} \langle v, \nabla f_{\beta_t}(x_t) \rangle$ 4:  $\gamma_t \leftarrow 2/(t+2)$ 5:  $x_{t+1} \leftarrow x_t + \gamma_t(v_t - x_t)$ 
  - 5:  $x_{t+1} \leftarrow x_t + \gamma_t (v_t x_t)$

### Lemma (Argyriou et al., 2014)

For all  $\beta \ge \beta' > 0$ ,  $f_{\beta} \leqslant f \leqslant f_{\beta} + \beta G_2^2/2$  and  $f_{\beta} \leqslant f_{\beta'} \leqslant f_{\beta} + (\beta - \beta')G_2^2/2$ .

### Theorem (Argyriou et al., 2014)

Let  $f:\mathbb{R}^n\to\mathbb{R}$  be a convex  $G_2\text{-Lipschitz}$  continuous function w.r.t. the  $\ell_2\text{-norm}.$  Then

$$f(x_t) - \min_{\mathcal{C}} f \leqslant \frac{4G_2D_2}{\sqrt{t+1}}$$

p	Assumption	Cardinality bound	
		Via Frank-Wolfe	Related work
[2,+∞[	_	$\mathcal{O}\left(\frac{pD_p^2}{\varepsilon^2}\right)$	$\mathcal{O}\left(\frac{pD_p^2}{\varepsilon^2}\right)$
	$x^* \in int(\mathcal{C})$	$\mathcal{O}\left(p\left(\frac{D_p}{r_p}\right)^2\ln\left(\frac{1}{\varepsilon}\right)\right)$	$\mathcal{O}\left(p\left(\frac{D_p}{r_p}\right)^2\ln\left(\frac{1}{\varepsilon}\right)\right)$
	${\mathcal C}$ strongly convex	$\mathcal{O}\left(\frac{\sqrt{p}D_p + p/\alpha_p}{\varepsilon}\right)$	- ,
]1,2[	-	$\mathcal{O}\left(rac{n^{(2-p)/p}D_2^2}{\varepsilon^2} ight)$	$\mathcal{O}\left(\frac{D_p^{p/(p-1)}}{p^{1/(p-1)}\varepsilon^{p/(p-1)}}\right)$
1	-	$\mathcal{O}\left(\frac{nD_2^2}{\varepsilon^2}\right)$	_
$+\infty$	_	$\mathcal{O}\left(\frac{D_2^2}{\varepsilon^2}\right)$	$\mathcal{O}\left(rac{\ln(n)D_{\infty}^2}{arepsilon^2} ight)$

## Conclusion

**Boosted Frank-Wolfe** 



Frank-Wolfe & adaptive gradients



Approximate Carathéodory



# References (1/4)

- A. Argyriou, M. Signoretto, and J. A. K. Suykens. Hybrid conditional gradient-smoothing algorithms with applications to sparse and low rank regularization. Chapman & Hall/CRC, 2014
- S. Barman. Approximating Nash equilibria and dense bipartite subgraphs via an approximate version of Carathéodory's theorem. *STOC*, 2015
- G. Braun, S. Pokutta, D. Tu, and S. Wright. Blended conditional gradients: the unconditioning of conditional gradients. *ICML*, 2019
- M. D. Canon and C. D. Cullum. A tight upper bound on the rate of convergence of Frank-Wolfe algorithm. SIAM J. Control, 1968
- C. W. Combettes and S. Pokutta. Blended matching pursuit. NeurIPS, 2019
- C. W. Combettes and S. Pokutta. Revisiting the approximate Carathéodory problem via the Frank-Wolfe algorithm. Minor revision at *Math. Program.*, 2021
- C. W. Combettes and S. Pokutta. Boosting Frank-Wolfe by chasing gradients. ICML, 2020
- C. W. Combettes and S. Pokutta. Complexity of linear minimization and projection on some sets. Minor revision at *Oper. Res. Lett.*, 2021
- C. W. Combettes, C. Spiegel, and S. Pokutta. Projection-free adaptive gradients for largescale optimization. *Submitted*, 2021
- L. Condat. Fast projection onto the simplex and the  $\ell_1$  ball. Math. Program., 2016

# References (2/4)

- V. F. Demyanov and A. M. Rubinov. *Approximate Methods in Optimization Problems*. Elsevier, 1970
- J. C. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. J. Mach. Learn. Res., 2011
- J. C. Dunn and S. Harshbarger. Conditional gradient algorithms with open loop step size rules. J. Math. Anal. Appl., 1978
- M. Frank and P. Wolfe. An algorithm for quadratic programming. Naval Res. Logist. Q., 1956
- D. Garber and O. Meshi. Linear-memory and decomposition-invariant linearly convergent conditional gradient algorithm for structured polytopes. *NIPS*, 2016
- E. Hazan and H. Luo. Variance-reduced and projection-free stochastic optimization. *ICML*, 2016
- J. Guélat and P. Marcotte. Some comments on Wolfe's 'away step'. Math. Program., 1986
- Y. Haugazeau. Sur les Inéquations Variationnelles et la Minimisation de Fonctionnelles Convexes. Thèse de doctorat, Univ. Paris, 1968
- M. Jaggi. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. ICML, 2013
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. ICLR, 2015
- S. Lacoste-Julien and M. Jaggi. On the global linear convergence of Frank-Wolfe optimization variants. *NIPS*, 2015

# References (3/4)

- G. Lan. The complexity of large-scale convex programming under a linear optimization oracle. *arXiv*, 2013
- G. Lan and Y. Zhou. Conditional gradient sliding for convex optimization. SIAM J. Optim., 2016
- E. S. Levitin and B. T. Polyak. Constrained minimization methods. USSR Comp. Math. Math. Phys., 1966
- F. Locatello, M. Tschannen, G. Rätsch, and M. Jaggi. Greedy algorithms for cone constrained optimization with convergence guarantees. *NIPS*, 2017
- H. B. McMahan and M. Streeter. Adaptive bound optimization for online convex optimization. *COLT*, 2010
- V. Mirrokni, R. Paes Leme, A. Vladu, and S. C.-W. Wong. Tight bounds for approximate Carathéodory and beyond. *ICML*, 2017
- J. J. Moreau. Proximité et dualité dans un espace hilbertien. Bull. Soc. Math. France, 1965
- G. Négiar, G. Dresdner, A. Y.-T. Tsai, L. El Ghaoui, F. Locatello, R. M. Freund, and F. Pedregosa. Stochastic Frank-Wolfe for constrained finite-sum minimization. ICML, 2020
- A. S. Nemirovsky and D. B. Yudin. Problem Complexity and Method Efficiency in Optimization. Wiley, 1983
- G. Pisier. Remarques sur un résultat non publié de B. Maurey. Ec. polytech., 1981
- S. J. Reddi, S. Kale, and S. Kumar. On the convergence of Adam and beyond. ICLR, 2018

- Z. Shen, C. Fang, P. Zhao, J. Huang, and H. Qian. Complexities in projection-free stochastic non-convex minimization. AISTATS, 2019
- P. Wolfe. Convergence theory in nonlinear programming. Integer and Nonlinear Programming. North-Holland, 1970
- J. Xie, Z. Shen, C. Zhang, H. Qian, and B. Wang. Efficient projection-free online methods with stochastic recursive gradient. AAAI, 2020
- A. Yurtsever, S. Sra, and V. Cevher. Conditional gradient methods via stochastic path-integrated differential estimator. ICML, 2019

Algorithm	Update $\tilde{\nabla} f(x_t)$
SFW	$rac{1}{b_t}\sum_{j=1}^{b_t} abla f_{i_j}(x_t)$
SVRF	$ abla f( ilde{x}_t) + rac{1}{b_t}\sum_{j=1}^{b_t} ( abla f_{i_j}(x_t) -  abla f_{i_j}( ilde{x}_t))$
SPIDER-FW	$ abla f( ilde{x}_t) + rac{1}{b_t} \sum_{i=1}^{b_t} ( abla f_{i_j}(x_t) -  abla f_{i_j}(x_{t-1}))$
ORGFW	$\frac{1}{b_t} \sum_{j=1}^{b_t} \nabla f_{i_j}(x_t) + (1 - \rho_t) \left( \tilde{\nabla} f(x_{t-1}) - \frac{1}{b_t} \sum_{j=1}^{b_t} \nabla f_{i_j}(x_{t-1}) \right)$
CSFW	$ ilde{ abla} f(x_{t-1}) + \sum_{i=1}^{^{\mathcal{D}_t}} \left( rac{1}{m} f_{i_j}'(\langle a_{i_j}, x_t  angle) - [lpha_{t-1}]_{i_j}  ight) a_{i_j}$
	and $[\alpha_t]_{i_j} \leftarrow \begin{cases} (1/m)f'_{i_j}(\langle a_{i_j},x_t\rangle) & \text{if } i_j \in \{i_1,\ldots,i_{b_t}\}\\ [\alpha_{t-1}]_{i_j} & \text{else} \end{cases}$

Algorithm	Update $\tilde{\nabla} f(x_t)$
SFW	$\frac{1}{b_t}\sum_{j=1}^{b_t}\nabla f_{i_j}(x_t)$
SVRF	$ abla f( ilde{x}_t) + rac{1}{b_t}\sum_{j=1}^{b_t} ( abla f_{i_j}(x_t) -  abla f_{i_j}( ilde{x}_t))$
SPIDER-FW	$ abla f( ilde{x}_t) + rac{1}{b_t} \sum_{j=1}^{b_t} ( abla f_{i_j}(x_t) -  abla f_{i_j}(x_{t-1}))$
ORGFW	$\frac{1}{b_t} \sum_{j=1}^{b_t} \nabla f_{i_j}(x_t) + (1 - \rho_t) \left( \tilde{\nabla} f(x_{t-1}) - \frac{1}{b_t} \sum_{j=1}^{b_t} \nabla f_{i_j}(x_{t-1}) \right)$
CSFW	$ ilde{ abla} f(x_{t-1}) + \sum_{j=1}^{b_t} \left( rac{1}{m} f_{i_j}'(\langle a_{i_j}, x_t  angle) - [lpha_{t-1}]_{i_j}  ight) a_{i_j}$
	and $[\alpha_t]_{i_j} \leftarrow \begin{cases} (1/m)f'_{i_j}(\langle a_{i_j}, x_t \rangle) & \text{if } i_j \in \{i_1, \dots, i_{b_t}\}\\ [\alpha_{t-1}]_{i_j} & \text{else} \end{cases}$

Algorithm	Update $\tilde{\nabla} f(x_t)$
SFW	$\frac{1}{b_t}\sum_{j=1}^{b_t}\nabla f_{i_j}(x_t)$
SVRF	$ abla f( ilde{x}_t) + rac{1}{b_t}\sum_{j=1}^{b_t} ( abla f_{i_j}(x_t) -  abla f_{i_j}( ilde{x}_t))$
SPIDER-FW	$ abla f( ilde{x}_t) + rac{1}{b_t}\sum_{j=1}^{b_t} ( abla f_{i_j}(x_t) -  abla f_{i_j}(x_{t-1}))$
ORGFW	$\frac{1}{b_t} \sum_{j=1}^{b_t} \nabla f_{i_j}(x_t) + (1 - \rho_t) \left( \tilde{\nabla} f(x_{t-1}) - \frac{1}{b_t} \sum_{j=1}^{b_t} \nabla f_{i_j}(x_{t-1}) \right)$
CSFW	$ ilde{ abla} f(x_{t-1}) + \sum_{j=1}^{b_t} \left( rac{1}{m} f_{ij}'(\langle a_{i_j}, x_t  angle) - [lpha_{t-1}]_{i_j}  ight) a_{i_j}$
	and $[\alpha_t]_{i_j} \leftarrow \begin{cases} (1/m)f'_{i_j}(\langle a_{i_j}, x_t \rangle) & \text{if } i_j \in \{i_1, \dots, i_{b_t}\} \\ [\alpha_{t-1}]_{i_j} & \text{else} \end{cases}$

Algorithm	Update $\tilde{\nabla} f(x_t)$
SFW	$\frac{1}{b_t}\sum_{j=1}^{b_t}\nabla f_{i_j}(x_t)$
SVRF	$ abla f( ilde{x}_t) + rac{1}{b_t}\sum_{j=1}^{b_t} ( abla f_{i_j}(x_t) -  abla f_{i_j}( ilde{x}_t))$
SPIDER-FW	$ abla f( ilde{x}_t) + rac{1}{b_t}\sum_{i=1}^{b_t} ( abla f_{i_j}( extsf{x}_t) -  abla f_{i_j}( extsf{x}_{t-1}))$
ORGFW	$\frac{1}{b_t} \sum_{j=1}^{b_t} \nabla f_{i_j}(x_t) + (1 - \rho_t) \left( \tilde{\nabla} f(x_{t-1}) - \frac{1}{b_t} \sum_{j=1}^{b_t} \nabla f_{i_j}(x_{t-1}) \right)$
CSFW	$ ilde{ abla} f(x_{t-1}) + \sum_{i=1}^{{}^{D_t}} \left( rac{1}{m} f_{ij}'(\langle a_{i_j}, x_t  angle) - [lpha_{t-1}]_{i_j}  ight) a_{i_j}$
	and $[\alpha_t]_{i_j} \leftarrow \begin{cases} (1/m)f'_{i_j}(\langle a_{i_j}, x_t \rangle) & \text{if } i_j \in \{i_1, \dots, i_{b_t}\} \\ [\alpha_{t-1}]_{i_j} & \text{else} \end{cases}$

Algorithm	Update $\tilde{\nabla} f(x_t)$
SFW	$rac{1}{b_t}\sum_{j=1}^{b_t} abla f_{i_j}(\mathbf{x}_t)$
SVRF	$ abla f( ilde{x}_t) + rac{1}{b_t}\sum_{j=1}^{b_t} ( abla f_{i_j}(x_t) -  abla f_{i_j}( ilde{x}_t))$
SPIDER-FW	$ abla f( ilde{x}_t) + rac{1}{b_t}\sum_{j=1}^{b_t} ( abla f_{i_j}(x_t) -  abla f_{i_j}(x_{t-1}))$
ORGFW	$\frac{1}{b_t} \sum_{j=1}^{b_t} \nabla f_{i_j}(x_t) + (1 - \rho_t) \left( \tilde{\nabla} f(x_{t-1}) - \frac{1}{b_t} \sum_{j=1}^{b_t} \nabla f_{i_j}(x_{t-1}) \right)$
CSFW	$ ilde{ abla} f(x_{t-1}) + \sum_{i=1}^{^{D_t}} \left( rac{1}{m} f_{ij}'(\langle \pmb{a}_{i_j}, x_t  angle) - [lpha_{t-1}]_{i_j}  ight) \pmb{a}_{i_j}$
	and $[\alpha_t]_{i_j} \leftarrow \begin{cases} (1/m)f'_{i_j}(\langle a_{i_j}, x_t \rangle) & \text{if } i_j \in \{i_1, \dots, i_{b_t}\} \\ [\alpha_{t-1}]_{i_j} & \text{else} \end{cases}$

Algorithm	Update $ ilde{ abla} f(x_t)$
SFW	$\frac{1}{b_t}\sum_{j=1}^{b_t}\nabla f_{i_j}(x_t)$
SVRF	$ abla f( ilde{x}_t) + rac{1}{b_t}\sum_{j=1}^{b_t} ( abla f_{i_j}(x_t) -  abla f_{i_j}( ilde{x}_t))$
SPIDER-FW	$ abla f( ilde{x}_t) + rac{1}{b_t}\sum_{j=1}^{b_t} ( abla f_{i_j}(x_t) -  abla f_{i_j}(x_{t-1}))$
ORGFW	$\frac{1}{b_t} \sum_{j=1}^{b_t} \nabla f_{i_j}(x_t) + (1 - \rho_t) \left( \tilde{\nabla} f(x_{t-1}) - \frac{1}{b_t} \sum_{j=1}^{b_t} \nabla f_{i_j}(x_{t-1}) \right)$
CSFW	$ ilde{ abla} f(\mathbf{x}_{t-1}) + \sum_{j=1}^{b_t} \left( rac{1}{m} f_{i_j}'(\langle \mathbf{a}_{i_j}, \mathbf{x}_t  angle) - [lpha_{t-1}]_{i_j}  ight) \mathbf{a}_{i_j}$
	and $[\alpha_t]_{i_j} \leftarrow \begin{cases} (1/m) f'_{i_j}(\langle a_{i_j}, x_t \rangle) & \text{if } i_j \in \{i_1, \dots, i_{b_t}\}\\ [\alpha_{t-1}]_{i_j} & \text{else} \end{cases}$

# A lower bound when $p \in [2, +\infty[$

### Theorem (Mirrokni et al., 2017)

Let  $p \in [2, +\infty[, H_n \in \mathbb{R}^{n \times n}$  be the Hadamard matrix of dimension n,  $C = \operatorname{conv}(H_n/n^{1/p})$  be the convex hull of the  $\ell_p$ -normalized columns of  $H_n$ , and  $x^* = e_1/n^{1/p} \in C$ . Then for all  $x \in C$ ,

$$\|x - x^*\|_p \leqslant \varepsilon \Rightarrow \operatorname{card}(x) \geqslant \frac{1}{\varepsilon^2 + 1/n}$$

# A lower bound when $p\in [2,+\infty[$

#### Theorem (Mirrokni et al., 2017)

Let  $p \in [2, +\infty[, H_n \in \mathbb{R}^{n \times n}$  be the Hadamard matrix of dimension n,  $C = \operatorname{conv}(H_n/n^{1/p})$  be the convex hull of the  $\ell_p$ -normalized columns of  $H_n$ , and  $x^* = e_1/n^{1/p} \in C$ . Then for all  $x \in C$ ,

$$\|x - x^*\|_p \leqslant \varepsilon \Rightarrow \operatorname{card}(x) \geqslant \frac{1}{\varepsilon^2 + 1/n}$$

•  $H_n \in \mathbb{R}^{n \times n}$  is a Hadamard matrix if  $H_n \in \{\pm 1\}^{n \times n}$  and  $H_n^\top H_n = nI_n$
#### Theorem (Mirrokni et al., 2017)

Let  $p \in [2, +\infty[, H_n \in \mathbb{R}^{n \times n}$  be the Hadamard matrix of dimension n,  $C = \operatorname{conv}(H_n/n^{1/p})$  be the convex hull of the  $\ell_p$ -normalized columns of  $H_n$ , and  $x^* = e_1/n^{1/p} \in C$ . Then for all  $x \in C$ ,

$$\|x - x^*\|_p \leqslant \varepsilon \Rightarrow \operatorname{card}(x) \geqslant rac{1}{\varepsilon^2 + 1/n}$$

- $H_n \in \mathbb{R}^{n \times n}$  is a Hadamard matrix if  $H_n \in \{\pm 1\}^{n \times n}$  and  $H_n^\top H_n = nI_n$
- Sylvester's construction:

$$H_{2n} = \begin{pmatrix} H_n & H_n \\ H_n & -H_n \end{pmatrix}$$

#### Theorem (Mirrokni et al., 2017)

Let  $p \in [2, +\infty[, H_n \in \mathbb{R}^{n \times n}$  be the Hadamard matrix of dimension n,  $C = \operatorname{conv}(H_n/n^{1/p})$  be the convex hull of the  $\ell_p$ -normalized columns of  $H_n$ , and  $x^* = e_1/n^{1/p} \in C$ . Then for all  $x \in C$ ,

$$\|x - x^*\|_p \leqslant \varepsilon \Rightarrow \operatorname{card}(x) \geqslant \frac{1}{\varepsilon^2 + 1/n}$$

- $H_n \in \mathbb{R}^{n \times n}$  is a Hadamard matrix if  $H_n \in \{\pm 1\}^{n \times n}$  and  $H_n^\top H_n = nI_n$
- Sylvester's construction:

$$H_{2n} = \begin{pmatrix} H_n & H_n \\ H_n & -H_n \end{pmatrix}$$

gives



FCFW almost matches the lower bound



- FCFW almost matches the lower bound
- There is no precise analysis of FCFW: the current analysis is transferred from that of AFW (Lacoste-Julien & Jaggi, 2015) and holds only for smooth strongly convex functions