

Splitting the Conditional Gradient Algorithm

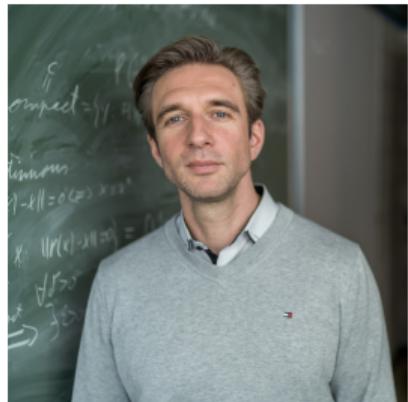
INFORMS Annual Meeting

Zev Woodstock

October 2023



Results are joint work with...



Sebastian Pokutta
ZIB & Technische
Universität Berlin



Interactive Optimization & Learning (IOL) Lab
iol.zib.de



Splitting the Conditional Gradient Algorithm

- 1.** Motivation: History of splitting and CG / “Frank-Wolfe” algorithms
- 2.** Algorithm design
- 3.** Convergence guarantees

1. Motivation: History of splitting and CG / "Frank-Wolfe" algorithms

Algorithms for one constraint

Classical problem setup

Given a smooth function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ and a nonempty compact convex set C ,

$$\text{minimize } f(x) \text{ subject to } x \in C. \quad (1)$$

Algorithms for one constraint

Classical problem setup

Given a smooth function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ and a nonempty **compact convex set C** ,

$$\text{minimize } f(x) \text{ subject to } x \in C. \quad (1)$$

Two iterative first-order algorithms for solving (1)

Projected gradient descent:

Conditional gradient:

Algorithms for one constraint

Classical problem setup

Given a smooth function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ and a nonempty **compact convex set C** ,

$$\text{minimize } f(x) \text{ subject to } x \in C. \quad (1)$$

Two iterative first-order algorithms for solving (1) differ in how $x \in C$ is enforced.

Projected gradient descent:

Conditional gradient:

Algorithms for one constraint

Classical problem setup

Given a smooth function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ and a nonempty **compact convex set C** ,

$$\text{minimize } f(x) \text{ subject to } x \in C. \quad (1)$$

Two iterative first-order algorithms for solving (1) differ in how $x \in C$ is enforced.

Projected gradient descent: Requires
the *projection onto C* , proj_C :

$$y \mapsto \arg \min_{x \in C} \|x - y\|^2 \quad (\text{PROJ})$$

Conditional gradient:

Algorithms for one constraint

Classical problem setup

Given a smooth function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ and a nonempty **compact convex set C** ,

$$\text{minimize } f(x) \text{ subject to } x \in C. \quad (1)$$

Two iterative first-order algorithms for solving (1) differ in how $x \in C$ is enforced.

Projected gradient descent: Requires the *projection onto C* , proj_C :

$$y \mapsto \arg \min_{x \in C} \|x - y\|^2 \quad (\text{PROJ})$$

Conditional gradient: Requires the *linear minimization oracle of C* , LMO_C :

$$y \mapsto p \in \arg \min_{x \in C} \langle y | x \rangle \quad (\text{LMO})$$

Algorithms for one constraint

Classical problem setup

Given a smooth function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ and a nonempty **compact convex set C** ,

$$\text{minimize } f(x) \text{ subject to } x \in C. \quad (1)$$

Two iterative first-order algorithms for solving (1) differ in how $x \in C$ is enforced.

Projected gradient descent: Requires the *projection onto C* , proj_C :

$$y \mapsto \arg \min_{x \in C} \|x - y\|^2 \quad (\text{PROJ})$$

Conditional gradient: Requires the *linear minimization oracle of C* , LMO_C :

$$y \mapsto p \in \arg \min_{x \in C} \langle y | x \rangle \quad (\text{LMO})$$

[Combettes/Pokutta, '21]: For many constraints, (PROJ) is **more expensive** than (LMO). (e.g., nuclear norm ball, ℓ_1 ball, probability simplex, Birkhoff polytope, general LP, ...)

What about multiple constraints?

Splitting problem setup

Given a smooth function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ and compact convex sets $(C_i)_{i \in I}$ ($I = \{1, \dots, m\}$),

$$\text{minimize } f(x) \text{ subject to } x \in \bigcap_{i \in I} C_i, \quad (*)$$

Applications: data science, matrix decomposition, quantum computing, combinatorial graph theory

What about multiple constraints?

Splitting problem setup

Given a smooth function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ and compact convex sets $(C_i)_{i \in I}$ ($I = \{1, \dots, m\}$),

$$\text{minimize } f(x) \text{ subject to } x \in \bigcap_{i \in I} C_i, \quad (*)$$

Applications: data science, matrix decomposition, quantum computing, combinatorial graph theory

Issue: Computing the projection or LMO for $\bigcap_{i \in I} C_i$ is prohibitively expensive.

What about multiple constraints?

Splitting problem setup

Given a smooth function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ and compact convex sets $(C_i)_{i \in I}$ ($I = \{1, \dots, m\}$),

$$\text{minimize } f(x) \text{ subject to } x \in \bigcap_{i \in I} C_i, \quad (*)$$

Applications: data science, matrix decomposition, quantum computing, combinatorial graph theory

Issue: Computing the projection or LMO for $\bigcap_{i \in I} C_i$ is prohibitively expensive.

Projection-based **splitting algorithms** (e.g., Forward-Backward, Douglas-Rachford, projective splitting, etc.), enforce constraints via projections onto the individual sets

Use $\text{proj}_{C_1}, \text{proj}_{C_2}, \dots$ instead of $\text{proj}(\bigcap_{i \in I} C_i)$

What about multiple constraints?

Splitting problem setup

Given a smooth function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ and compact convex sets $(C_i)_{i \in I}$ ($I = \{1, \dots, m\}$),

$$\text{minimize } f(x) \text{ subject to } x \in \bigcap_{i \in I} C_i, \quad (*)$$

Applications: data science, matrix decomposition, quantum computing, combinatorial graph theory

Issue: Computing the projection or LMO for $\bigcap_{i \in I} C_i$ is prohibitively expensive.

Projection-based **splitting algorithms** (e.g., Forward-Backward, Douglas-Rachford, projective splitting, etc.), enforce constraints via projections onto the individual sets

Use $\text{proj}_{C_1}, \text{proj}_{C_2}, \dots$ instead of $\text{proj}(\bigcap_{i \in I} C_i)$

Simpler tools → previously intractable problems become solvable on a larger scale.

1. Motivation: History of splitting and CG / "Frank-Wolfe" algorithms

What if projections are too expensive?

LMO-based *splitting algorithms*, enforce constraints via LMOs for the individual sets

Use $\text{LMO}_{C_1}, \text{LMO}_{C_2}, \dots$ instead of $\text{LMO}_{(\bigcap_{i \in I} C_i)}$



$\text{LMO}_{\bigcap_{i \in I} C_i}$



$\text{LMO}_{C_1}, \dots, \text{LMO}_{C_m}$

1. Motivation: History of splitting and CG / "Frank-Wolfe" algorithms

What if projections are too expensive?

LMO-based *splitting algorithms*, enforce constraints via LMOs for the individual sets

Use $\text{LMO}_{C_1}, \text{LMO}_{C_2}, \dots$ instead of $\text{LMO}_{(\bigcap_{i \in I} C_i)}$

Relatively little has been done in this field.



$\text{LMO}_{\bigcap_{i \in I} C_i}$



$\text{LMO}_{C_1}, \dots, \text{LMO}_{C_m}$

1. Motivation: History of splitting and CG / "Frank-Wolfe" algorithms

What if projections are too expensive?

LMO-based *splitting algorithms*, enforce constraints via LMOs for the individual sets

Use $\text{LMO}_{C_1}, \text{LMO}_{C_2}, \dots$ instead of $\text{LMO}_{(\bigcap_{i \in I} C_i)}$

Relatively little has been done in this field.

- Unlike projections, LMOs are discontinuous.
- “CTRL+F / Replace proj with LMO” fails.



$\text{LMO}_{\bigcap_{i \in I} C_i}$



$\text{LMO}_{C_1}, \dots, \text{LMO}_{C_m}$

1. Motivation: History of splitting and CG / "Frank-Wolfe" algorithms

What if projections are too expensive?

LMO-based *splitting algorithms*, enforce constraints via LMOs for the individual sets

Use $\text{LMO}_{C_1}, \text{LMO}_{C_2}, \dots$ instead of $\text{LMO}_{(\bigcap_{i \in I} C_i)}$

Relatively little has been done in this field.

- Unlike projections, LMOs are discontinuous.
- "CTRL+F / Replace proj with LMO" fails.
- "State-of-the-art" relies on inexact prox-based algorithms (mostly Augmented Lagrangians).



$\text{LMO}_{\bigcap_{i \in I} C_i}$

$\text{LMO}_{C_1}, \dots, \text{LMO}_{C_m}$

1. Motivation: History of splitting and CG / "Frank-Wolfe" algorithms

Previous work

"Use a CG subroutine to approximate a projection" \Rightarrow high iteration complexity
[He/Harchaoui, '15], [Liu et al., '19] [Millan et al., '21], [Kolmogorov/Pock, '21]

1. Motivation: History of splitting and CG / "Frank-Wolfe" algorithms

Previous work

"Use a CG subroutine to approximate a projection" \Rightarrow high iteration complexity
[He/Harchaoui, '15], [Liu et al., '19] [Millan et al., '21], [Kolmogorov/Pock, '21]

Currently, lowest iteration complexity is $\mathcal{O}(m)$: one LMO per set.

1. Motivation: History of splitting and CG / "Frank-Wolfe" algorithms

Previous work

"Use a CG subroutine to approximate a projection" \Rightarrow high iteration complexity
[He/Harchaoui, '15], [Liu et al., '19] [Millan et al., '21], [Kolmogorov/Pock, '21]

Currently, lowest iteration complexity is $\mathcal{O}(m)$: one LMO per set.

	$m = 2$	$m > 2$	f convex	f nonconvex	Analysis
[Pedregosa et al., '20]	✗	✗	✓	✓	CG
[Braun et al., '22]	✓	✗	✗ ($f = 0$)	✗	CG

Previous work

"Use a CG subroutine to approximate a projection" \Rightarrow high iteration complexity
 [He/Harchaoui, '15], [Liu et al., '19] [Millan et al., '21], [Kolmogorov/Pock, '21]

Currently, lowest iteration complexity is $\mathcal{O}(m)$: one LMO per set.

	$m = 2$	$m > 2$	f convex	f nonconvex	Analysis
[Pedregosa et al., '20]	✗	✗	✓	✓	CG
[Braun et al., '22]	✓	✗	✗($f = 0$)	✗	CG
[Gidel et al. '18]	✓	✗	✓	✗	AL+CG

|

(✓)- requires additional structure on $(C_i)_{i \in I}$

1. Motivation: History of splitting and CG / "Frank-Wolfe" algorithms

Previous work

"Use a CG subroutine to approximate a projection" \Rightarrow high iteration complexity
 [He/Harchaoui, '15], [Liu et al., '19] [Millan et al., '21], [Kolmogorov/Pock, '21]

Currently, lowest iteration complexity is $\mathcal{O}(m)$: one LMO per set.

	$m = 2$	$m > 2$	f convex	f nonconvex	Analysis
[Pedregosa et al., '20]	✗	✗	✓	✓	CG
[Braun et al., '22]	✓	✗	✗ ($f = 0$)	✗	CG
[Gidel et al. '18]	✓	✗	✓	✗	AL+CG
[Yurtsever et al. '19], [Silvetti-Falls et al. '20]	✓	✓	✓	✗	AL+CG
[Lan et al., '21]	(✓)	(✓)	✓	✗	CG

(✓)- requires additional structure on $(C_i)_{i \in I}$

1. Motivation: History of splitting and CG / "Frank-Wolfe" algorithms

Previous work

"Use a CG subroutine to approximate a projection" \Rightarrow high iteration complexity
 [He/Harchaoui, '15], [Liu et al., '19] [Millan et al., '21], [Kolmogorov/Pock, '21]

Currently, lowest iteration complexity is $\mathcal{O}(m)$: one LMO per set.

	$m = 2$	$m > 2$	f convex	f nonconvex	Analysis
[Pedregosa et al., '20]	✗	✗	✓	✓	CG
[Braun et al., '22]	✓	✗	✗ ($f = 0$)	✗	CG
[Gidel et al. '18]	✓	✗	✓	✗	AL+CG
[Yurtsever et al. '19], [Silvetti-Falls et al. '20]	✓	✓	✓	✗	AL+CG
[Lan et al., '21]	(✓)	(✓)	✓	✗	CG
[ZW/Pokutta '23]	✓	✓	✓	✓	CG

(✓)- requires additional structure on $(C_i)_{i \in I}$

Splitting the Conditional Gradient Algorithm

- 1.** Motivation: History of splitting and CG / “Frank-Wolfe” algorithms
- 2.** Algorithm design
- 3.** Convergence guarantees

Tools from the projection literature

Product space construction

- $\{\omega_i\}_{i \in I} \subset]0, 1]$, $\sum_{i \in I} \omega_i = 1$ (e.g., $\omega_i \equiv 1/m$)
- $\mathcal{H} = \mathbb{R}^n$ and $\mathcal{H} = \times_{i \in I} \mathcal{H}$, with inner product $\sum_{i \in I} \omega_i \langle \cdot | \cdot \rangle$
- *Diagonal subspace of \mathcal{H} :* $D = \{(x, \dots, x) \mid x \in \mathcal{H}\}$

Projecting onto D amounts to computing an average

$$\text{proj}_D x = A^* A = A^* \sum_{i \in I} \omega_i x^i.$$

Tools from the projection literature

Product space construction

- $\{\omega_i\}_{i \in I} \subset]0, 1]$, $\sum_{i \in I} \omega_i = 1$ (e.g., $\omega_i \equiv 1/m$)
- $\mathcal{H} = \mathbb{R}^n$ and $\mathcal{H} = \times_{i \in I} \mathcal{H}_i$, with inner product $\sum_{i \in I} \omega_i \langle \cdot | \cdot \rangle$
- *Diagonal subspace of \mathcal{H} :* $\mathcal{D} = \{(x, \dots, x) \mid x \in \mathcal{H}\}$
- *Block-averaging operator* and its adjoint:

$$A: \mathcal{H} \rightarrow \mathcal{H}: (\mathbf{x}^1, \dots, \mathbf{x}^m) \mapsto \sum_{i \in I} \omega_i \mathbf{x}^i \quad A^*: \mathbf{x} \mapsto (x, \dots, x).$$

Projecting onto \mathcal{D} amounts to computing an average

$$\text{proj}_{\mathcal{D}} \mathbf{x} = A^* A = A^* \sum_{i \in I} \omega_i \mathbf{x}^i.$$

Tools from the projection literature

Product space construction

- $\{\omega_i\}_{i \in I} \subset]0, 1]$, $\sum_{i \in I} \omega_i = 1$ (e.g., $\omega_i \equiv 1/m$)
- $\mathcal{H} = \mathbb{R}^n$ and $\mathcal{H} = \times_{i \in I} \mathcal{H}$, with inner product $\sum_{i \in I} \omega_i \langle \cdot | \cdot \rangle$
- *Diagonal subspace of \mathcal{H} :* $\mathbf{D} = \{(x, \dots, x) \mid x \in \mathcal{H}\}$
- *Block-averaging operator* and its adjoint:

$$A: \mathcal{H} \rightarrow \mathcal{H}: (\mathbf{x}^1, \dots, \mathbf{x}^m) \mapsto \sum_{i \in I} \omega_i \mathbf{x}^i \quad A^*: \mathbf{x} \mapsto (x, \dots, x).$$

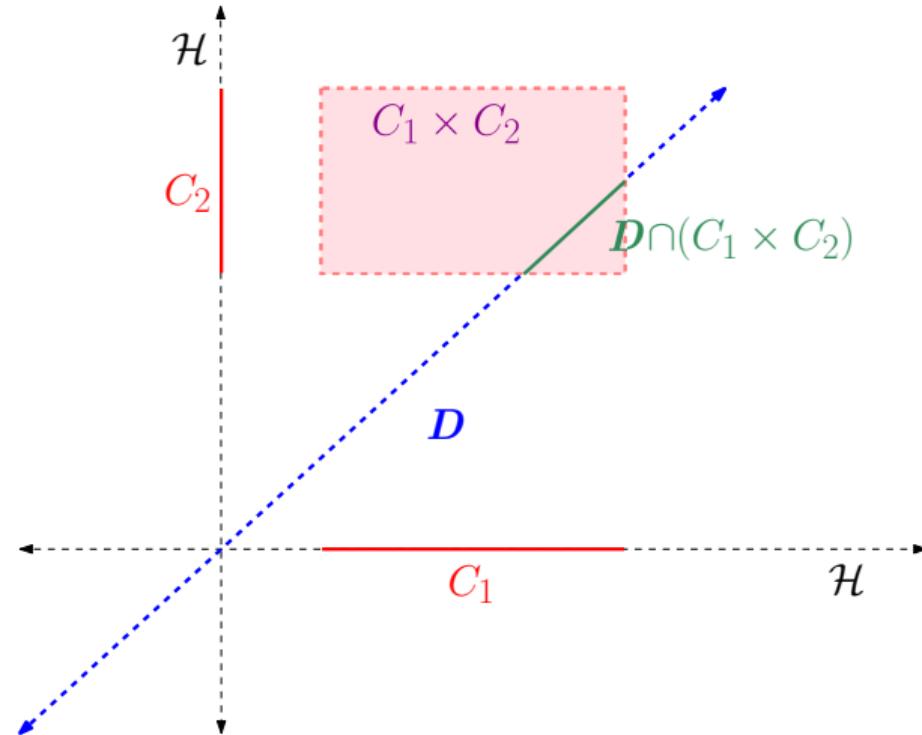
Projecting onto \mathbf{D} amounts to computing an average

$$\text{proj}_{\mathbf{D}} \mathbf{x} = A^* A = A^* \sum_{i \in I} \omega_i \mathbf{x}^i.$$

Tools from the projection literature

Product space construction

- $\mathcal{H} = \times_{i \in I} \mathcal{H}$
- $D = \{(x, \dots, x) \mid x \in \mathcal{H}\} \subset \mathcal{H}$



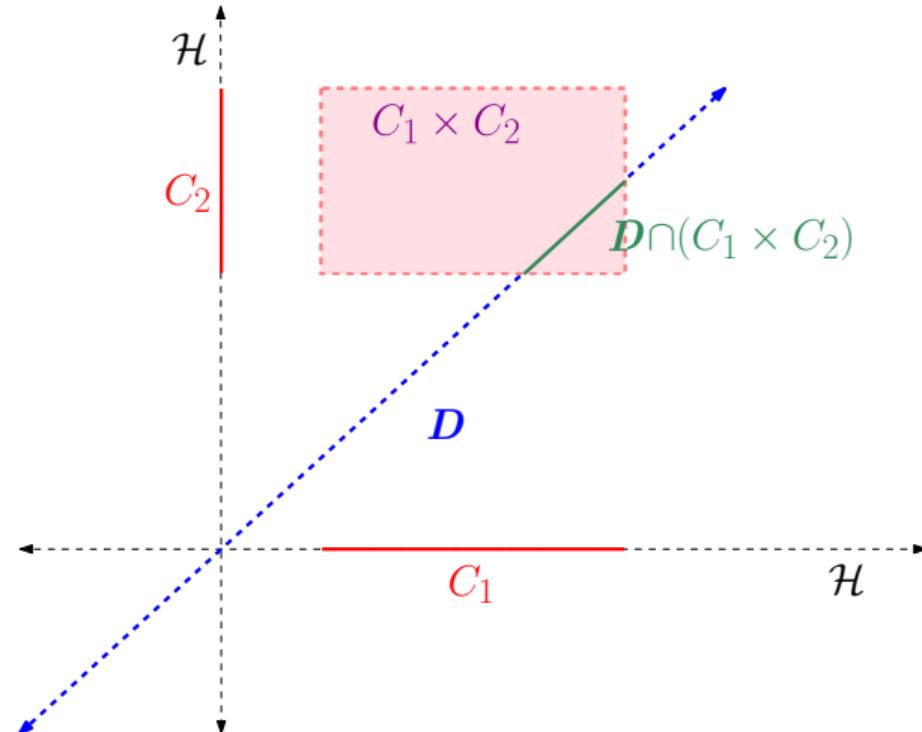
Tools from the projection literature

Product space construction

- $\mathcal{H} = \times_{i \in I} \mathcal{H}$
- $D = \{(x, \dots, x) \mid x \in \mathcal{H}\} \subset \mathcal{H}$

Proposition (Reformulation of $\bigcap_{i \in I} C_i$)

$x \in D \cap \times_{i \in I} C_i$ if and only if
 $x = (x, \dots, x)$ and $x \in \bigcap_{i \in I} C_i$



2. Algorithm design

Product space relaxation

$$\text{minimize } f(x) \quad \text{subject to} \quad \textcolor{teal}{x} \in \bigcap_{i \in I} C_i, \quad (*)$$

$$\underset{x \in \times_{i \in I} C_i}{\text{minimize}} \quad \underbrace{f(Ax) + \lambda_t \text{dist}_D^2(x)}_{F_{\lambda_t}(x)}. \quad (*)$$

Product space relaxation

$$\text{minimize } f(x) \quad \text{subject to} \quad \textcolor{teal}{x} \in \bigcap_{i \in I} C_i, \tag{*}$$

admits the equivalent reformulation (via the $0-\infty$ indicator function ι_D)

$$\underset{x \in D \cap \bigcap_{i \in I} C_i}{\text{minimize}} \quad f(Ax) = \underset{x \in \bigcap_{i \in I} C_i}{\text{minimize}} \quad f(Ax) + \iota_D(x).$$

$$\underset{x \in \bigcap_{i \in I} C_i}{\text{minimize}} \quad \underbrace{f(Ax) + \lambda_t \text{dist}_D^2(x)}_{F_{\lambda_t}(x)}. \quad (*)$$

Product space relaxation

$$\text{minimize } f(x) \quad \text{subject to} \quad \textcolor{teal}{x} \in \bigcap_{i \in I} C_i, \tag{*}$$

admits the equivalent reformulation (via the $0-\infty$ indicator function ι_D)

$$\underset{x \in D \cap \bigcap_{i \in I} C_i}{\text{minimize}} \quad f(Ax) = \underset{x \in \bigcap_{i \in I} C_i}{\text{minimize}} \quad f(Ax) + \iota_D(x).$$

Relaxation (for $\lambda_t \geq 0$)

$$\underset{x \in \bigcap_{i \in I} C_i}{\text{minimize}} \quad \underbrace{f(Ax) + \lambda_t \text{dist}_D^2(x)}_{F_{\lambda_t}(x)}. \tag{*}$$

Product space relaxation

$$\text{minimize } f(x) \quad \text{subject to} \quad \textcolor{teal}{x} \in \bigcap_{i \in I} C_i, \tag{*}$$

admits the equivalent reformulation (via the $0-\infty$ indicator function ι_D)

$$\underset{x \in D \cap \bigcap_{i \in I} C_i}{\text{minimize}} \quad f(Ax) = \underset{x \in \bigcap_{i \in I} C_i}{\text{minimize}} \quad f(Ax) + \iota_D(x).$$

Relaxation (for $\lambda_t \geq 0$)

→ Relaxation is tractable with vanilla CG!

$$\underset{x \in \bigcap_{i \in I} C_i}{\text{minimize}} \quad \underbrace{f(Ax) + \lambda_t \text{dist}_D^2(x)}_{F_{\lambda_t}(x)}. \tag{*}$$

$$\nabla F_{\lambda_t}(x) = A^* \nabla f(Ax) + \lambda_t(x - \text{proj}_D x)$$

$$\text{LMO}_{\bigcap_{i \in I} C_i}(x) = (\text{LMO}_{C_1}(x^1), \dots, \text{LMO}_{C_m}(x^m))$$

Product space relaxation

$$\text{minimize } f(x) \quad \text{subject to} \quad \textcolor{teal}{x} \in \bigcap_{i \in I} C_i, \tag{*}$$

admits the equivalent reformulation (via the $0-\infty$ indicator function ι_D)

$$\underset{x \in D \cap \bigcap_{i \in I} C_i}{\text{minimize}} \quad f(Ax) = \underset{x \in \bigcap_{i \in I} C_i}{\text{minimize}} \quad f(Ax) + \iota_D(x).$$

Relaxation (for $\lambda_t \geq 0$)

→ Relaxation is tractable with vanilla CG!

$$\underset{x \in \bigcap_{i \in I} C_i}{\text{minimize}} \quad \underbrace{f(Ax) + \lambda_t \text{dist}_D^2(x)}_{F_{\lambda_t}(x)}. \tag{*}$$

$$\nabla F_{\lambda_t}(x) = A^* \nabla f(Ax) + \lambda_t(x - \text{proj}_D x)$$

$$\text{LMO}_{\bigcap_{i \in I} C_i}(x) = (\text{LMO}_{C_1}(x^1), \dots, \text{LMO}_{C_m}(x^m))$$

Pseudocode:

Product space relaxation

$$\text{minimize } f(x) \quad \text{subject to} \quad \textcolor{teal}{x} \in \bigcap_{i \in I} C_i, \tag{*}$$

admits the equivalent reformulation (via the $0-\infty$ indicator function ι_D)

$$\underset{x \in D \cap \bigcap_{i \in I} C_i}{\text{minimize}} \quad f(Ax) = \underset{x \in \bigcap_{i \in I} C_i}{\text{minimize}} \quad f(Ax) + \iota_D(x).$$

Relaxation (for $\lambda_t \geq 0$)

→ Relaxation is tractable with vanilla CG!

$$\underset{x \in \bigcap_{i \in I} C_i}{\text{minimize}} \quad \underbrace{f(Ax) + \lambda_t \text{dist}_D^2(x)}_{F_{\lambda_t}(x)}. \tag{*}$$

$$\nabla F_{\lambda_t}(x) = A^* \nabla f(Ax) + \lambda_t(x - \text{proj}_D x)$$

$$\text{LMO}_{\bigcap_{i \in I} C_i}(x) = (\text{LMO}_{C_1}(x^1), \dots, \text{LMO}_{C_m}(x^m))$$

Pseudocode:(A) Perform one CG step on (*);

Product space relaxation

$$\text{minimize } f(x) \quad \text{subject to} \quad \textcolor{teal}{x} \in \bigcap_{i \in I} C_i, \quad (*)$$

admits the equivalent reformulation (via the $0-\infty$ indicator function ι_D)

$$\underset{x \in D \cap \bigcap_{i \in I} C_i}{\text{minimize}} \quad f(Ax) = \underset{x \in \bigcap_{i \in I} C_i}{\text{minimize}} \quad f(Ax) + \iota_D(x).$$

Relaxation (for $\lambda_t \geq 0$)

→ Relaxation is tractable with vanilla CG!

$$\underset{x \in \bigcap_{i \in I} C_i}{\text{minimize}} \quad \underbrace{f(Ax) + \lambda_t \text{dist}_D^2(x)}_{F_{\lambda_t}(x)}. \quad (*)$$

$$\nabla F_{\lambda_t}(x) = A^* \nabla f(Ax) + \lambda_t(x - \text{proj}_D x)$$

$$\text{LMO}_{\bigcap_{i \in I} C_i}(x) = (\text{LMO}_{C_1}(x^1), \dots, \text{LMO}_{C_m}(x^m))$$

Pseudocode: (A) Perform one CG step on (*); (B) Update λ_t ;

Product space relaxation

$$\text{minimize } f(x) \quad \text{subject to} \quad \textcolor{teal}{x} \in \bigcap_{i \in I} C_i, \quad (*)$$

admits the equivalent reformulation (via the $0-\infty$ indicator function ι_D)

$$\underset{x \in D \cap \bigcap_{i \in I} C_i}{\text{minimize}} \quad f(Ax) = \underset{x \in \bigcap_{i \in I} C_i}{\text{minimize}} \quad f(Ax) + \iota_D(x).$$

Relaxation (for $\lambda_t \geq 0$)

→ Relaxation is tractable with vanilla CG!

$$\underset{x \in \bigcap_{i \in I} C_i}{\text{minimize}} \quad \underbrace{f(Ax) + \lambda_t \text{dist}_D^2(x)}_{F_{\lambda_t}(x)}. \quad (*)$$

$$\nabla F_{\lambda_t}(x) = A^* \nabla f(Ax) + \lambda_t(x - \text{proj}_D x)$$

$$\text{LMO}_{\bigcap_{i \in I} C_i}(x) = (\text{LMO}_{C_1}(x^1), \dots, \text{LMO}_{C_m}(x^m))$$

Pseudocode: (A) Perform one CG step on (*); (B) Update λ_t ; (C) $t \leftarrow t + 1$.

The algorithm

Split Conditional Gradient (SCG) Algorithm

Require: Point $x_0 \in \sum_{i \in I} \omega_i C_i$, smooth function f , weights $\{\omega_i\}_{i \in I} \subset]0, 1]$ such that $\sum_{i \in I} \omega_i = 1$

```
1: for  $t = 0, 1$  to ... do
2:   Choose penalty parameter  $\lambda_t \in ]0, +\infty[$ 
3:   Choose step size  $\gamma_t \in ]0, 1]$ 
4:    $g_t \leftarrow \nabla f(x_t)$ 
5:   for  $i = 1$  to  $m$  do
6:      $v_t^i \leftarrow \text{LMO}_i(g_t + \lambda_t(x_t^i - x_t))$ 
7:      $x_{t+1}^i \leftarrow x_t^i + \gamma_t(v_t^i - x_t^i)$ 
8:   end for
9:    $x_{t+1} \leftarrow \sum_{i \in I} \omega_i x_{t+1}^i$ 
10:  end for
```

Practical advantages:

- Uses individual LMOs
- m LMO calls per iteration.
- Line 9: speeds up feasibility.

The algorithm

Split Conditional Gradient (SCG) Algorithm

Require: Point $x_0 \in \sum_{i \in I} \omega_i C_i$, smooth function f , weights $\{\omega_i\}_{i \in I} \subset]0, 1]$ such that $\sum_{i \in I} \omega_i = 1$

```
1: for  $t = 0, 1$  to ... do
2:   Choose penalty parameter  $\lambda_t \in ]0, +\infty[$ 
3:   Choose step size  $\gamma_t \in ]0, 1]$ 
4:    $g_t \leftarrow \nabla f(x_t)$ 
5:   for  $i = 1$  to  $m$  do
6:      $v_t^i \leftarrow \text{LMO}_i(g_t + \lambda_t(x_t^i - x_t))$ 
7:      $x_{t+1}^i \leftarrow x_t^i + \gamma_t(v_t^i - x_t^i)$ 
8:   end for
9:    $x_{t+1} \leftarrow \sum_{i \in I} \omega_i x_{t+1}^i$ 
10:  end for
```

Practical advantages:

- Uses individual LMOs
- m LMO calls per iteration.
- Line 9: speeds up feasibility.

Question:

- Does it actually solve (\star) ?

TL;DR: Yes.

$\gamma_t = \mathcal{O}(1/\sqrt{t})$ and $\lambda_t = \mathcal{O}(\ln t)$ work.

Why does averaging help?

$$\mathbf{x} \in \mathcal{D} \cap \bigtimes_{i \in I} C_i \Rightarrow A\mathbf{x} \in \bigcap_{i \in I} C_i,$$

so a feasible average is easier to satisfy than a feasible component!

Why does averaging help?

$$\mathbf{x} \in D \cap \bigtimes_{i \in I} C_i \Rightarrow A\mathbf{x} \in \bigcap_{i \in I} C_i,$$

so a feasible average is easier to satisfy than a feasible component!

Proposition

$A\mathbf{x}_t \in \bigcap_{i \in I} C_i$ if and only if
 $\text{proj}_D(\mathbf{x}) \in \bigtimes_{i \in I} C_i.$

Why does averaging help?

$$\mathbf{x} \in \mathcal{D} \cap \bigtimes_{i \in I} C_i \Rightarrow A\mathbf{x} \in \bigcap_{i \in I} C_i,$$

so a feasible average is easier to satisfy than a feasible component!

Proposition

$A\mathbf{x}_t \in \bigcap_{i \in I} C_i$ if and only if
 $\text{proj}_{\mathcal{D}}(\mathbf{x}) \in \bigtimes_{i \in I} C_i.$

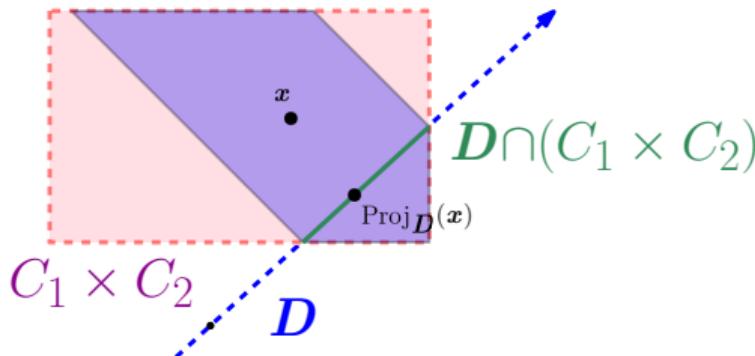


Figure: Darker shaded region $\{\mathbf{x} \in \mathcal{H} \mid A\mathbf{x} \in \bigcap_{i \in I} C_i\}$ contains the segment $\mathcal{D} \cap \bigtimes_{i \in I} C_i.$

Splitting the Conditional Gradient Algorithm

- 1.** Motivation: History of splitting and CG / “Frank-Wolfe” algorithms
- 2.** Algorithm design
- 3.** Convergence guarantees

Convergence of the subproblems

Proposition (Convergence of relaxed problems)

Let $(\lambda_t)_{t \in \mathbb{N}} \rightarrow +\infty$. For every $t \in \mathbb{N}$, set $F_t = f \circ A + \lambda_t \text{dist}_D^2 / 2 + \iota_{X_{i \in I} C_i}$. Then

1. F_t converges pointwise to $f \circ A + \iota_{D \cap X_{i \in I} C_i}$.
2. F_t converges epigraphically to $f \circ A + \iota_{D \cap X_{i \in I} C_i}$.
3. ∂F_t converges graphically to $\partial(f \circ A + \iota_{D \cap X_{i \in I} C_i})$.

where epigraphical and graphical convergence are in, e.g., [Rockafellar/Wets, '09].

Proposition (Convergence of optimal values for $\lambda_t \nearrow +\infty$)

$$\lim_{t \rightarrow +\infty} \left(\inf_{x \in X_{i \in I} C_i} F_{\lambda_t}(x) \right) \rightarrow \inf_{x \in X_{i \in I} C_i} \left(\lim_{t \rightarrow +\infty} F_{\lambda_t}(x) \right) = \inf_{x \in \bigcap_{i \in I} C_i} f(x).$$

Convex case

Theorem (Convex convergence)

Let f be convex and L_f -smooth, let $(C_i)_{i \in I}$ be nonempty compact convex subsets of \mathcal{H} with diameters $\{R_i\}_{i \in I} \subset [0, +\infty[$ such that $\bigcap_{i \in I} C_i \neq \emptyset$, and for every $\lambda \geq 0$, set $F_\lambda: \mathbf{x} \mapsto f(A\mathbf{x}) + \frac{\lambda}{2} \text{dist}_D^2(\mathbf{x})$. Let $\lambda_0 > 0$ and $\lambda_{t+1} = \lambda_t + (\sqrt{t} + 2)^{-2}$ and $\gamma_t = 2/(\sqrt{t} + 2)$. Then

$$0 \leq F_{\lambda_t}(\mathbf{x}_t) - F_{\lambda_t}(\mathbf{x}_t^*) \leq \mathcal{O}\left(\frac{\ln t}{\sqrt{t}}\right)$$

In particular,

1. $F_{\lambda_t}(\mathbf{x}_t) \rightarrow \inf_{x \in \bigcap_{i \in I} C_i} f(x)$ and $\text{dist}_D(\mathbf{x}_t) \rightarrow 0$.
2. Every accumulation point \mathbf{x}_∞ of $(\mathbf{x}_t)_{t \in \mathbb{N}}$ produces a solution $A\mathbf{x}_\infty \in \bigcap_{i \in I} C_i$ such that $f(A\mathbf{x}_\infty) = \inf_{x \in \bigcap_{i \in I} C_i} f(x)$.

Convex case

Theorem (Convex convergence)

Let f be convex and L_f -smooth, let $(C_i)_{i \in I}$ be nonempty compact convex subsets of \mathcal{H} with diameters $\{R_i\}_{i \in I} \subset [0, +\infty[$ such that $\bigcap_{i \in I} C_i \neq \emptyset$, and for every $\lambda \geq 0$, set $F_\lambda: \mathbf{x} \mapsto f(A\mathbf{x}) + \frac{\lambda}{2} \text{dist}_D^2(\mathbf{x})$. Let $\lambda_0 > 0$ and $\lambda_{t+1} = \lambda_t + (\sqrt{t} + 2)^{-2}$ and $\gamma_t = 2/(\sqrt{t} + 2)$. Then

$$0 \leq F_{\lambda_t}(\mathbf{x}_t) - F_{\lambda_t}(\mathbf{x}_t^*) \leq \mathcal{O}\left(\frac{\ln t}{\sqrt{t}}\right)$$

In particular,

1. $F_{\lambda_t}(\mathbf{x}_t) \rightarrow \inf_{x \in \bigcap_{i \in I} C_i} f(x)$ and $\text{dist}_D(\mathbf{x}_t) \rightarrow 0$.
2. Every accumulation point \mathbf{x}_∞ of $(\mathbf{x}_t)_{t \in \mathbb{N}}$ produces a solution $A\mathbf{x}_\infty \in \bigcap_{i \in I} C_i$ such that $f(A\mathbf{x}_\infty) = \inf_{x \in \bigcap_{i \in I} C_i} f(x)$.

We believe this rate can be improved!

Nonconvex case

Theorem (Nonconvex convergence)

Let f be L_f -smooth, let $(C_i)_{i \in I}$ be nonempty compact convex subsets of \mathcal{H} with diameters $\{R_i\}_{i \in I} \subset [0, +\infty[$ such that $\bigcap_{i \in I} C_i \neq \emptyset$, and for every $\lambda \geq 0$, set $F_\lambda: \mathbf{x} \mapsto f(A\mathbf{x}) + \frac{\lambda}{2} \text{dist}_D^2(\mathbf{x})$. Let $\lambda_t = \sum_{k=0}^{t-1} 1/(k+1)$ and $\gamma_t = 1/\sqrt{t}$. Then,

$$0 \leq \frac{1}{t} \sum_{k=0}^{t-1} \langle \nabla F_{\lambda_k}(\mathbf{x}_k) \mid \mathbf{x}_k - \mathbf{v}_k \rangle \leq \mathcal{O}\left(\frac{\ln t}{\sqrt{t}} + \frac{1}{\sqrt{t}}\right).$$

In particular, there exists a subsequence $(t_k)_{k \in \mathbb{N}}$ such that

1. $(\langle \nabla F_{\lambda_{t_k}}(\mathbf{x}_{t_k}) \mid \mathbf{x}_{t_k} - \mathbf{v}_{t_k} \rangle)_{k \in \mathbb{N}} \rightarrow 0$ and $\text{dist}_D(\mathbf{x}_{t_k}) \rightarrow 0$.
2. Furthermore, every accumulation point \mathbf{x}_∞ of $(\mathbf{x}_{t_k})_{k \in \mathbb{N}}$ yields a stationary point $A\mathbf{x}_\infty \in \bigcap_{i \in I} C_i$ of the problem (\star) .

3. Convergence guarantees

Best-known rates / Future work

	$m = 2$	$m > 2$	f convex	f nonconvex
[Pedregosa et al., '20]	✗	✗	$\mathcal{O}\left(\frac{1}{t}\right)$	$\mathcal{O}\left(\frac{1}{\sqrt{t}}\right)$
[Gidel et al. '18]	✓	✗	$\mathcal{O}\left(\frac{1}{t}\right)$	✗
[Yurtsever et al. '19], [Lan et al., '21] (✓)	✓	✓	$\mathcal{O}\left(\frac{1}{\sqrt{t}}\right)$	✗
[ZW/Pokutta '23]	✓	✓	✓	$\mathcal{O}\left(\frac{\ln t}{\sqrt{t}}\right)$

3. Convergence guarantees

Best-known rates / Future work

	$m = 2$	$m > 2$	f convex	f nonconvex	
[Pedregosa et al., '20]	✗	✗	$\mathcal{O}\left(\frac{1}{t}\right)$	$\mathcal{O}\left(\frac{1}{\sqrt{t}}\right)$	←
[Gidel et al. '18]	✓	✗	$\mathcal{O}\left(\frac{1}{t}\right)$	✗	
[Yurtsever et al. '19], [Lan et al., '21] (✓)	✓	✓	$\mathcal{O}\left(\frac{1}{\sqrt{t}}\right)$	✗	
[ZW/Pokutta '23]	✓	✓	✓	$\mathcal{O}\left(\frac{\ln t}{\sqrt{t}}\right)$	←

Usually, there is a quadratic speed-up from nonconvex and convex rates.

Best-known rates / Future work

	$m = 2$	$m > 2$	f convex	f nonconvex	
[Pedregosa et al., '20]	✗	✗	$\mathcal{O}\left(\frac{1}{t}\right)$	$\mathcal{O}\left(\frac{1}{\sqrt{t}}\right)$	
[Gidel et al. '18]	✓	✗	$\mathcal{O}\left(\frac{1}{t}\right)$	✗	←
[Yurtsever et al. '19], [Lan et al., '21] (✓)	✓	✓	$\mathcal{O}\left(\frac{1}{\sqrt{t}}\right)$	✗	←
[ZW/Pokutta '23]	✓	✓	✓	$\mathcal{O}\left(\frac{\ln t}{\sqrt{t}}\right)$	

Usually, there is a quadratic speed-up from nonconvex and convex rates.

Is the gap between $m = 2$ and $m > 2$ actually necessary?

Thank you for your attention!

References

-  G. Braun, S. Pokutta, and R. Weismantel, Alternating linear minimization: revisiting von Neumann's alternating projections
preprint, arXiv: 2212.02933
-  R. Díaz Millán and O. P. Ferreira and L. F. Prudente, Alternating conditional gradient method for convex feasibility problems
Comput. Optim. Appl., vol. 80, pp. 245–269, 2021
-  G. Gidel, F. Pedregosa, and S. Lacoste-Julien, Frank-Wolfe splitting via augmented Lagrangian Method
Proc. AISTATS, pp. 1456–1465, 2018.
-  N. He and Z. Harchaoui, Semi-proximal mirror-prox for nonsmooth composite minimization
Proc. NeurIPS, vol. 28, 2015
-  V. Kolmogorov and T. Pock, One-sided Frank-Wolfe algorithms for saddle problems
Proc. ICML, PMLR, vol. 139, pp. 5665–5675, 2021
-  G. Lan, E. Romeijn, and Z. Zhou, Conditional gradient methods for convex optimization with general affine and nonlinear constraints
SIAM J. Optim., vol. 31, no. 3, pp. 2307–2339, 2021.

References

-  Y-F. Liu, X. Liu, and S. Ma, On the nonergodic convergence rate of an inexact augmented Lagrangian framework for composite convex programming
Math. Oper. Res., vol. 44, no. 2 pp. 632–650, 2019
-  F. Pedregosa, G. Negiar, A. Askari, and M. Jaggi, Linearly convergent Frank-Wolfe with backtracking line-search
Proc. AISTATS, pp. 1–10, 2020.
-  R. T. Rockafellar, and R. J-B Wets, Variational Analysis
Springer, 2009
-  A. Silveti-Falls, C. Molinari, and J. Fadili, Linearly convergent Frank-Wolfe with backtracking line-search
SIAM J. Optim., vol. 30, no. 4, pp. 2687–2725, 2020.
-  A. Yurtsever, O. Fercoq, and V. Cevher, A conditional-gradient-based augmented Lagrangian framework
Proc. ICML, pp. 7272–7281, 2019.