

# A quest for theoretically-sound optimization

Zev Woodstock

Technische Universität Berlin & Zuse Institute Berlin

February 2024



# “How did I get here?” -David Byrne



## “How did I get here?” -David Byrne



Lise

6.05 – 8 petaflop/second  
(roughly 75 – 100 IBM Watsons)



# Theoretically-sound optimization

1. Motivation
2. Background: Theory vs practice
3. Proximity operators: Algorithmic bells and whistles
4. Splitting FW: What if the “usual” tools fail us?
5. More adventuring



# What is optimization?

## Optimization in a nutshell ( $\mathcal{H} = \mathbb{R}^n$ or any real Hilbert space)

- **Objective function**  $f: \mathcal{H} \rightarrow \mathbb{R} \cup \{+\infty\}$ .  
e.g., data fidelity in ML, energy efficiency, profit, statistical error, ...
- An “optimal”  $x \in \mathcal{H}$  makes  $f(x)$  the smallest or largest  
e.g., **minimize** error, **maximize** efficiency

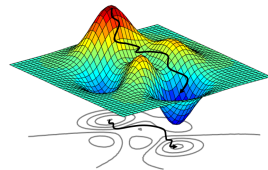


image: towardsdatascience.com

$$\underset{x \in \mathcal{H}}{\text{minimize}} f(x)$$

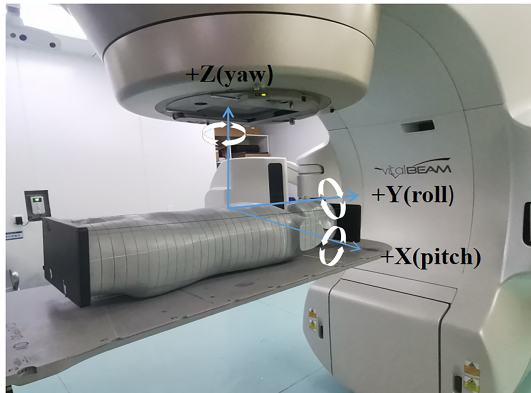
**Constraint set(s)**  $\mathcal{C} \subset \mathcal{H}$

e.g.,  $\mathbb{R}_+^N$ ,  $\mathbb{S}_+^N$ , hypercube, solution set of an  
inverse problem, ...

$$\iota_{\mathcal{C}}(x) = \begin{cases} 0 & \text{if } x \in \mathcal{C} \\ +\infty & \text{otherwise.} \end{cases}$$

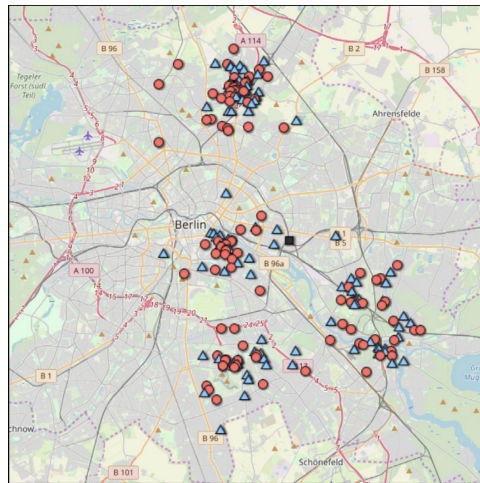
$$\underset{x \in \mathcal{C}}{\text{minimize}} \tilde{f}(x) = \underset{x \in \mathcal{H}}{\text{minimize}} \underbrace{\tilde{f}(x) + \iota_{\mathcal{C}}(x)}_f$$

## Modeling via optimization



[Torelli et al., *Med. Phys.*, 2023]

image: [Fu et al., *Tech. Cancer Res. Treatment*, 2023]



[Sartori & Buriol, *Comput. Oper. Res.*, 2020]

# Modeling via optimization



0	2	15	0	0	11	10	0	0	0	0	9	9	0	0	0
0	0	0	4	60	157	236	255	255	177	95	61	32	0	0	29
0	10	16	119	238	255	244	245	243	250	249	255	222	103	10	0
0	14	170	255	255	244	254	255	253	245	255	249	253	251	124	1
2	98	255	228	255	251	254	211	141	116	122	215	251	238	255	49
13	217	243	255	155	33	226	52	2	0	10	13	232	255	255	36
16	229	252	254	49	12	0	0	7	7	0	70	237	252	235	62
6	141	245	255	212	25	11	9	3	0	115	236	243	255	137	0
0	87	252	250	248	215	60	0	1	121	252	255	248	144	6	0
0	13	113	255	255	245	255	182	181	248	252	242	208	36	0	19
1	0	5	117	251	255	241	255	247	255	241	162	17	0	7	0
0	0	0	4	58	251	255	246	254	253	255	120	11	0	1	0
0	0	0	4	17	255	255	255	248	252	255	244	255	182	10	0
0	22	206	252	246	251	241	100	24	113	255	245	255	194	9	0
0	111	255	242	255	158	24	0	0	6	39	255	232	230	56	0
0	218	251	250	137	7	11	0	0	0	2	62	255	250	125	3
0	173	255	255	101	9	20	0	13	3	13	182	251	245	61	0
0	107	251	241	255	230	98	55	19	118	217	248	253	255	52	4
0	18	146	250	255	247	255	255	255	249	255	240	255	129	0	5
0	23	113	215	255	250	248	255	255	248	248	118	14	12	0	0
0	0	6	1	0	52	153	233	255	252	147	37	0	0	4	1
0	0	5	5	0	0	0	0	14	1	0	6	6	0	0	0

0	2	15	0	0	11	10	0	0	0	0	9	9	0	0	0
0	0	0	4	60	157	236	255	255	177	95	61	32	0	0	29
0	10	16	119	238	255	244	245	243	250	249	255	222	103	10	0
0	14	170	255	255	244	254	255	253	245	255	249	253	251	124	1
2	98	255	228	255	251	254	211	141	116	122	215	251	238	255	49
13	217	243	255	155	33	226	52	2	0	10	13	232	255	255	36
16	229	252	254	49	12	0	0	7	7	0	70	237	252	235	62
6	141	245	255	212	25	11	9	3	0	115	236	243	255	137	0
0	87	252	250	248	215	60	0	1	121	252	255	248	144	6	0
0	13	113	255	255	245	255	182	181	248	252	242	208	36	0	19
1	0	5	117	251	255	241	255	247	255	241	162	17	0	7	0
0	0	0	4	58	251	255	246	254	253	255	120	11	0	1	0
0	0	0	4	17	255	255	255	248	252	255	244	255	182	10	0
0	22	206	252	246	251	241	100	24	113	255	245	255	194	9	0
0	111	255	242	255	158	24	0	0	6	39	255	232	230	56	0
0	218	251	250	137	7	11	0	0	0	2	62	255	250	125	3
0	173	255	255	101	9	20	0	13	3	13	182	251	245	61	0
0	107	251	241	255	230	98	55	19	118	217	248	253	255	52	4
0	18	146	250	255	247	255	255	255	249	255	240	255	129	0	5
0	23	113	215	255	250	248	255	255	248	248	118	14	12	0	0
0	0	6	1	0	52	153	233	255	252	147	37	0	0	4	1
0	0	5	5	0	0	0	0	14	1	0	6	6	0	0	0

image: [towardsdatascience.com](https://towardsdatascience.com)

# Modeling via optimization

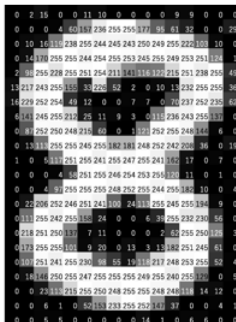
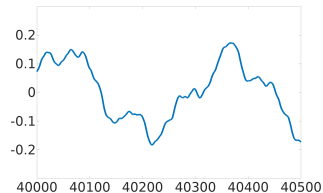
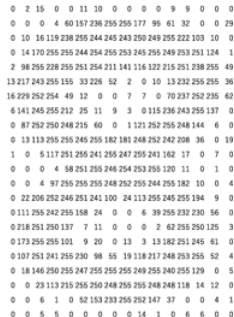


image: towardsdatascience.com



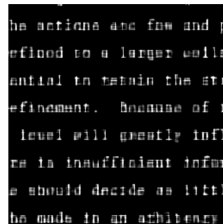
# Modeling via optimization

he actions are few and p  
efined to a larger colle  
ential to retain the str  
efinement. Because of t  
level will greatly infl  
re is insufficient infor  
e should decide as littl  
be made in an arbitrary

Original



Observation



Recovery

# Modeling via optimization

he actions are few and p  
efined to a larger colle  
ential to retain the str  
efinement. Because of t  
level will greatly infl  
re is insufficient infor  
e should decide as littl  
be made in an arbitrary

Original



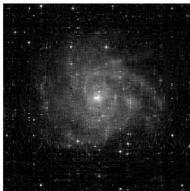
Observation

he actions are few and p  
efined to a larger colle  
ential to retain the str  
efinement. Because of t  
level will greatly infl  
re is insufficient infor  
e should decide as littl  
be made in an arbitrary

Recovery



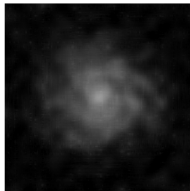
Original



Observation



Recovered stars



Recovered galaxy

[Combettes & ZW., *SIAM J. Imaging Sci.*, 2022]

# Some fundamental questions



- What are the roadblocks to *provably* solving optimization problems?
  - Nonconvexity, nonsmoothness, and bears – oh my!

# Some fundamental questions



- What are the roadblocks to *provably* solving optimization problems?
  - Nonconvexity, nonsmoothness, and bears – oh my!
- What *theoretically-sound* algorithms exist?



# Some fundamental questions



- What are the roadblocks to *provably* solving optimization problems?
  - Nonconvexity, nonsmoothness, and bears – oh my!
- What *theoretically-sound* algorithms exist, and can we do better?
  - Splitting, Parallelization, Extrapolation, Asynchronous computation

# Theoretically-sound optimization

1. Motivation
2. Background: Theory vs practice
3. Proximity operators: Algorithmic bells and whistles
4. Splitting FW: What if the “usual” tools fail us?
5. More adventuring

## Our setting

$\mathcal{H}$  is a real Hilbert space with inner product  $\langle \cdot | \cdot \rangle$ ,  
e.g.,  $\mathbb{R}^n$  with the dot product  $\langle x | y \rangle = x^T y$ .

## Our setting

$\mathcal{H}$  is a real Hilbert space with inner product  $\langle \cdot | \cdot \rangle$ ,  
e.g.,  $\mathbb{R}^n$  with the dot product  $\langle x | y \rangle = x^T y$ .

$\text{Argmin } f = \{x \in \mathcal{H} \mid f(x) = \inf f(\mathcal{H})\}$  is the set of minimizers of  $f$ .

## Our setting

$\mathcal{H}$  is a real Hilbert space with inner product  $\langle \cdot | \cdot \rangle$ ,  
e.g.,  $\mathbb{R}^n$  with the dot product  $\langle x | y \rangle = x^T y$ .

$\text{Argmin } f = \{x \in \mathcal{H} \mid f(x) = \inf f(\mathcal{H})\}$  is the set of minimizers of  $f$ .

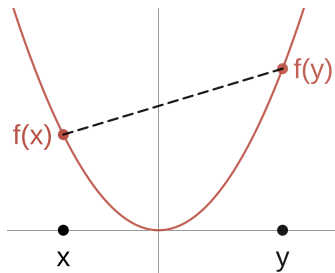
$f$  is **convex** if, for all  $x, y \in \mathcal{H}$  and  $\alpha \in (0, 1)$ ,  $f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$

## Our setting

$\mathcal{H}$  is a real Hilbert space with inner product  $\langle \cdot | \cdot \rangle$ ,  
e.g.,  $\mathbb{R}^n$  with the dot product  $\langle x | y \rangle = x^T y$ .

$\text{Argmin } f = \{x \in \mathcal{H} \mid f(x) = \inf f(\mathcal{H})\}$  is the set of minimizers of  $f$ .

$f$  is **convex** if, for all  $x, y \in \mathcal{H}$  and  $\alpha \in (0, 1)$ ,  $f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$

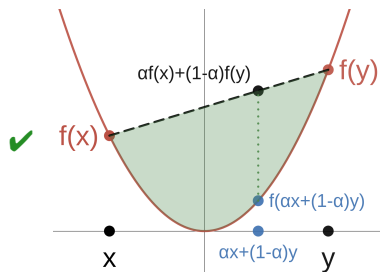


## Our setting

$\mathcal{H}$  is a real Hilbert space with inner product  $\langle \cdot | \cdot \rangle$ ,  
e.g.,  $\mathbb{R}^n$  with the dot product  $\langle x | y \rangle = x^T y$ .

$\text{Argmin } f = \{x \in \mathcal{H} \mid f(x) = \inf f(\mathcal{H})\}$  is the set of minimizers of  $f$ .

$f$  is **convex** if, for all  $x, y \in \mathcal{H}$  and  $\alpha \in (0, 1)$ ,  $f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$

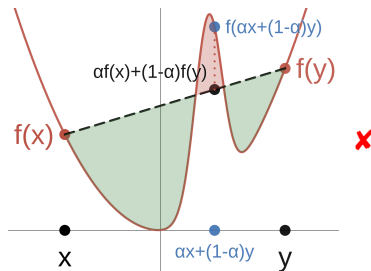
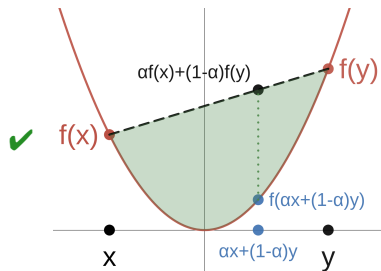


## Our setting

$\mathcal{H}$  is a real Hilbert space with inner product  $\langle \cdot | \cdot \rangle$ ,  
e.g.,  $\mathbb{R}^n$  with the dot product  $\langle x | y \rangle = x^T y$ .

$\text{Argmin } f = \{x \in \mathcal{H} \mid f(x) = \inf f(\mathcal{H})\}$  is the set of minimizers of  $f$ .

$f$  is **convex** if, for all  $x, y \in \mathcal{H}$  and  $\alpha \in (0, 1)$ ,  $f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$



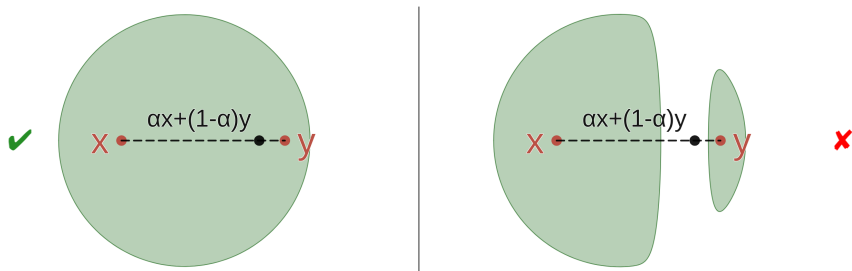


## Our setting

$\mathcal{H}$  is a real Hilbert space with inner product  $\langle \cdot | \cdot \rangle$ ,  
e.g.,  $\mathbb{R}^n$  with the dot product  $\langle x | y \rangle = x^T y$ .

$\text{Argmin } f = \{x \in \mathcal{H} \mid f(x) = \inf f(\mathcal{H})\}$  is the set of minimizers of  $f$ .

$C$  is **convex** if, for all  $x, y \in C$  and  $\alpha \in (0, 1)$ ,  $\alpha x + (1 - \alpha)y \in C$ .



## Our setting

$\mathcal{H}$  is a real Hilbert space with inner product  $\langle \cdot | \cdot \rangle$ ,  
e.g.,  $\mathbb{R}^n$  with the dot product  $\langle x | y \rangle = x^T y$ .

$\text{Argmin } f = \{x \in \mathcal{H} \mid f(x) = \inf f(\mathcal{H})\}$  is the set of minimizers of  $f$ .

$f$  is **convex** if, for all  $x, y \in \mathcal{H}$  and  $\alpha \in (0, 1)$ ,  $f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$

Let  $\Gamma_0(\mathcal{H}) = \{f : \mathcal{H} \rightarrow \mathbb{R} \cup \{+\infty\} \mid f \text{ is convex, lower-semicontinuous, and proper}\}$ ,

## Our setting

$\mathcal{H}$  is a real Hilbert space with inner product  $\langle \cdot | \cdot \rangle$ ,  
e.g.,  $\mathbb{R}^n$  with the dot product  $\langle x | y \rangle = x^T y$ .

$\text{Argmin } f = \{x \in \mathcal{H} \mid f(x) = \inf f(\mathcal{H})\}$  is the set of minimizers of  $f$ .

$f$  is **convex** if, for all  $x, y \in \mathcal{H}$  and  $\alpha \in (0, 1)$ ,  $f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$

Let  $\Gamma_0(\mathcal{H}) = \{f : \mathcal{H} \rightarrow \mathbb{R} \cup \{+\infty\} \mid f \text{ is convex, lower-semicontinuous, and proper}\}$ ,  
e.g.,  $e^x$ ,  $-\ln(x)$ ,  $\|\cdot\|^2$ , ReLU, Hinge loss,  $\|Ax + b\|$ ,  $\|\cdot\|_1$ ,  $\iota_C$  ( $C$  convex and closed),

## Our setting

$\mathcal{H}$  is a real Hilbert space with inner product  $\langle \cdot | \cdot \rangle$ ,  
e.g.,  $\mathbb{R}^n$  with the dot product  $\langle x | y \rangle = x^T y$ .

$\text{Argmin } f = \{x \in \mathcal{H} \mid f(x) = \inf f(\mathcal{H})\}$  is the set of minimizers of  $f$ .

$f$  is **convex** if, for all  $x, y \in \mathcal{H}$  and  $\alpha \in (0, 1)$ ,  $f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$

Let  $\Gamma_0(\mathcal{H}) = \{f : \mathcal{H} \rightarrow \mathbb{R} \cup \{+\infty\} \mid f \text{ is convex, lower-semicontinuous, and proper}\}$ ,

e.g.,  $e^x$ ,  $-\ln(x)$ ,  $\|\cdot\|^2$ , ReLU, Hinge loss,  $\|Ax + b\|$ ,  $\|\cdot\|_1$ ,  $\iota_C$  ( $C$  convex and closed),  
 $\sup\{f_i \mid i \in I\}$ , affine composition, positive linear combinations, ...

## Our setting

$\mathcal{H}$  is a real Hilbert space with inner product  $\langle \cdot | \cdot \rangle$ ,  
e.g.,  $\mathbb{R}^n$  with the dot product  $\langle x | y \rangle = x^T y$ .

$\text{Argmin } f = \{x \in \mathcal{H} \mid f(x) = \inf f(\mathcal{H})\}$  is the set of minimizers of  $f$ .

$f$  is **convex** if, for all  $x, y \in \mathcal{H}$  and  $\alpha \in (0, 1)$ ,  $f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$

Let  $\Gamma_0(\mathcal{H}) = \{f : \mathcal{H} \rightarrow \mathbb{R} \cup \{+\infty\} \mid f \text{ is convex, lower-semicontinuous, and proper}\}$ ,

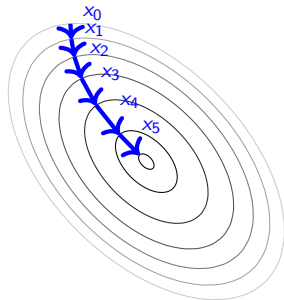
e.g.,  $e^x$ ,  $-\ln(x)$ ,  $\|\cdot\|^2$ , ReLU, Hinge loss,  $\|Ax + b\|$ ,  $\|\cdot\|_1$ ,  $\iota_C$  ( $C$  convex and closed),  
 $\sup\{f_i \mid i \in I\}$ , affine composition, positive linear combinations, ...

not  $\|\mathcal{N}(x) - d\|$  for multilayer neural networks

# “Traditioooooon”

-Tevye, Fiddler on the Roof

$f$  is  **$L$ -smooth** ( $L \geq 0$ ) if it is differentiable and  $\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$ .



# “Traditioooooon”

-Tevye, Fiddler on the Roof

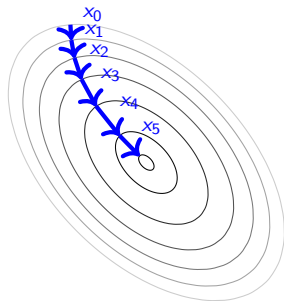
$f$  is  **$L$ -smooth** ( $L \geq 0$ ) if it is differentiable and  $\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$ .

## Gradient Descent

Let  $f: \mathcal{H} \rightarrow \mathbb{R}$  be  $L$ -smooth and suppose  $\text{Argmin } f \neq \emptyset$ . Let  $x_0 \in \mathcal{H}$ ,  $\varepsilon > 0$  and for every  $n \in \mathbb{N}$ , set

$$x_{n+1} = x_n - \lambda_n \nabla f(x_n), \quad \text{where } \lambda_n \in \left[ \varepsilon, \frac{2}{L} - \varepsilon \right] \quad (\text{GD})$$

If  $f \in \Gamma_0(\mathcal{H})$ , then  $(x_n)_{n \in \mathbb{N}}$  converges to a minimizer of  $f$ .



## Foe #1: Non-convexity

$f$  is  **$L$ -smooth** ( $L \geq 0$ ) if it is differentiable and  $\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$ .

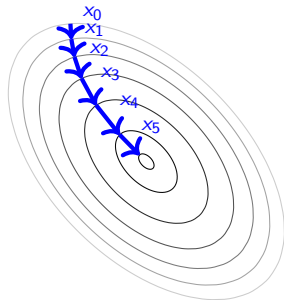
### Gradient Descent

Let  $f: \mathcal{H} \rightarrow \mathbb{R}$  be  $L$ -smooth and suppose  $\text{Argmin } f \neq \emptyset$ . Let  $x_0 \in \mathcal{H}$ ,  $\varepsilon > 0$  and for every  $n \in \mathbb{N}$ , set

$$x_{n+1} = x_n - \lambda_n \nabla f(x_n), \quad \text{where } \lambda_n \in \left[ \varepsilon, \frac{2}{L} - \varepsilon \right] \quad (\text{GD})$$

If  $f \notin \Gamma_0(\mathcal{H})$ , then  $(x_n)_{n \in \mathbb{N}}$  converges to a **stationary point**.  
 $(\nabla f(x^*) = 0)$

For  $x_0$  **sufficiently close** to a minimizer,  $(x_n)_{n \in \mathbb{N}}$  converges to one.





# “Traditioooooon”

-Tevye, Fiddler on the Roof

$f$  is  **$L$ -smooth** ( $L \geq 0$ ) if it is differentiable and  $\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$ .

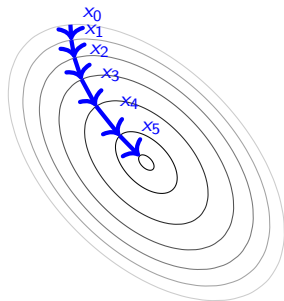
$$f = \sum_{1 \leq i \leq m} f_i$$

## Gradient Descent

Let  $f: \mathcal{H} \rightarrow \mathbb{R}$  be  $L$ -smooth and suppose  $\text{Argmin } f \neq \emptyset$ . Let  $x_0 \in \mathcal{H}$ ,  $\varepsilon > 0$  and for every  $n \in \mathbb{N}$ , set

$$x_{n+1} = x_n - \lambda_n \nabla f(x_n), \quad \text{where } \lambda_n \in \left[ \varepsilon, \frac{2}{L} - \varepsilon \right] \quad (\text{GD})$$

If  $f \in \Gamma_0(\mathcal{H})$ , then  $(x_n)_{n \in \mathbb{N}}$  converges to a minimizer of  $f$ .



# “Traditioooooon”

-Tevye, Fiddler on the Roof

$f$  is  **$L$ -smooth** ( $L \geq 0$ ) if it is differentiable and  $\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$ .

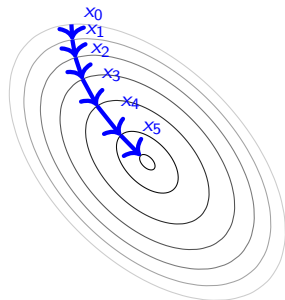
$$f = \sum_{1 \leq i \leq m} f_i$$

## Stochastic Gradient Descent (one variant)

Let  $f: \mathcal{H} \rightarrow \mathbb{R}$  be  $L$ -smooth and suppose  $\text{Argmin } f \neq \emptyset$ . Let  $x_0 \in \mathcal{H}$  and for every  $n \in \mathbb{N}$ , set

$$x_{n+1} = x_n - \frac{1}{n+1} \nabla f_{i_n}(x_n), \quad \text{where } i_n \sim U(\{1, \dots, m\}) \quad (\text{SGD})$$

If  $f \in \Gamma_0(\mathcal{H})$ , then  $\mathbb{E}[f(x_n)]$  converges to  $\inf_{x \in \mathcal{H}} f(x)$ .



# Why can't we take the eagles to Mordor?

(A reasonable question to ask, if we didn't read the books)

A common paradigm:

1. Define an **objective function**
2. Optimize with an efficient algorithm, e.g., SGD with **algorithmic differentiation** (AD)

# Why can't we take the eagles to Mordor?

(A reasonable question to ask, if we didn't read the books)

A common paradigm:

1. Define an **objective function**
2. Optimize with an efficient algorithm, e.g., SGD with **algorithmic differentiation** (AD)

[Pontil et al., *Numer. Algorithms*, 2019]

Training a sparse linear binary classifier

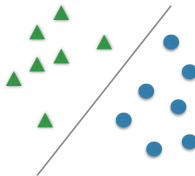


image: adeveloperdiary.com

# Why can't we take the eagles to Mordor?

(A reasonable question to ask, if we didn't read the books)

A common paradigm:

1. Define an **objective function**
2. Optimize with an efficient algorithm, e.g., SGD with **algorithmic differentiation** (AD)

[Pontil et al., *Numer. Algorithms*, 2019]

Training a sparse linear binary classifier

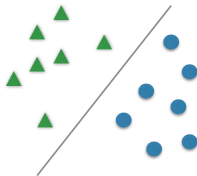


image: adeveloperdiary.com

$$\begin{aligned} \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad & \sum_{i \in I_1} \max\{0, 1 - \langle x \mid a_i \rangle\} + \\ & \sum_{i \in I_2} \max\{0, 1 + \langle x \mid a_i \rangle\} + \lambda \|x\|_1 \end{aligned}$$

## Foe #2: Non-differentiability

A common paradigm:

1. Define an **objective function**
2. Optimize with an efficient algorithm, e.g., SGD with **algorithmic differentiation** (AD)

**Issue:** For many **objective functions**, a **gradient** does not exist.

[Pontil et al., *Numer. Algorithms*, 2019]

Training a sparse linear binary classifier

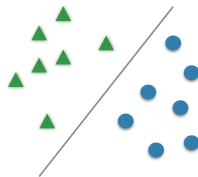


image: adeveloperdiary.com

$$\begin{aligned} \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad & \sum_{i \in I_1} \max\{0, 1 - \langle x \mid a_i \rangle\} + \\ & \sum_{i \in I_2} \max\{0, 1 + \langle x \mid a_i \rangle\} + \lambda \|x\|_1 \end{aligned}$$

## Foe #2: Non-differentiability

A common paradigm:

1. Define an **objective function**
2. Optimize with an efficient algorithm, e.g., SGD with **algorithmic differentiation** (AD)

**Issue:** For many **objective functions**, a **gradient** does not exist.

Engineers\*:



image: ripleys.com

\*:Some ~~mathematicians at heart~~ exceptions exist

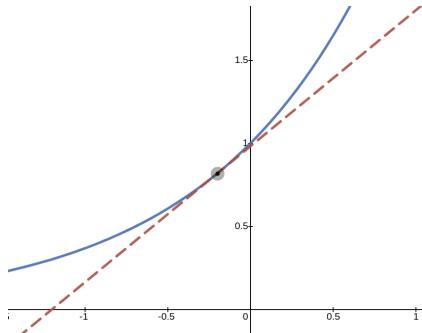
## How do we solve $\nabla f = 0$ when $\nabla f$ doesn't exist?



## How do we solve $\nabla f = 0$ when $\nabla f$ doesn't exist?

If  $f \in \Gamma_0(\mathcal{H})$  is differentiable at  $x \in \mathcal{H}$ , then

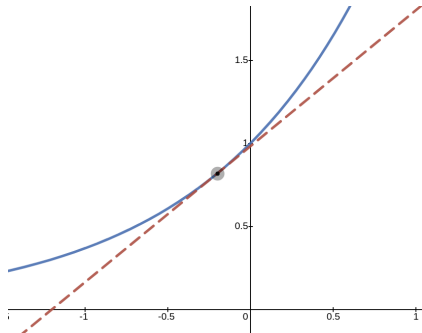
$$(\forall y \in \mathcal{H}) \quad \langle y - x \mid \nabla f(x) \rangle + f(x) \leq f(y).$$



## How do we solve $\nabla f = 0$ when $\nabla f$ doesn't exist?

A **subgradient**  $g \in \mathcal{H}$  of  $f: \mathcal{H} \rightarrow ]-\infty, +\infty]$  at  $x \in \mathcal{H}$  satisfies

$$(\forall y \in \mathcal{H}) \quad \langle y - x \mid g \rangle + f(x) \leq f(y).$$

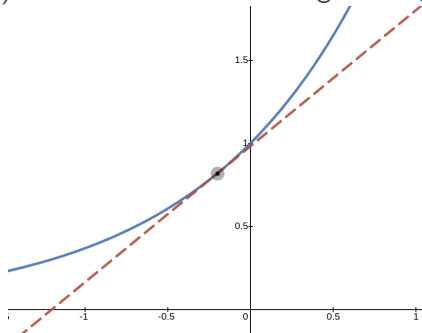


## How do we solve $\nabla f = 0$ when $\nabla f$ doesn't exist?

A **subgradient**  $g \in \mathcal{H}$  of  $f: \mathcal{H} \rightarrow ]-\infty, +\infty]$  at  $x \in \mathcal{H}$  satisfies

$$(\forall y \in \mathcal{H}) \quad \langle y - x \mid g \rangle + f(x) \leq f(y).$$

The **subdifferential**  $\partial f(x) \subset \mathcal{H}$  is the set containing **all subgradients** of  $f$  at  $x$ .



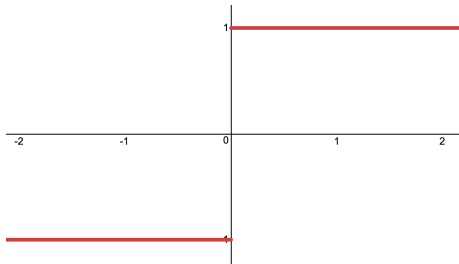
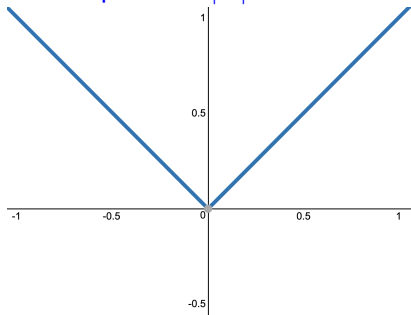
## How do we solve $\nabla f = 0$ when $\nabla f$ doesn't exist?

A **subgradient**  $g \in \mathcal{H}$  of  $f: \mathcal{H} \rightarrow ]-\infty, +\infty]$  at  $x \in \mathcal{H}$  satisfies

$$(\forall y \in \mathcal{H}) \quad \langle y - x \mid g \rangle + f(x) \leq f(y).$$

The **subdifferential**  $\partial f(x) \subset \mathcal{H}$  is the set containing **all subgradients** of  $f$  at  $x$ .

**Example:**  $f = |\cdot|$ : What do we do at zero?



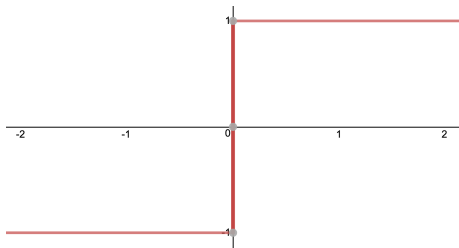
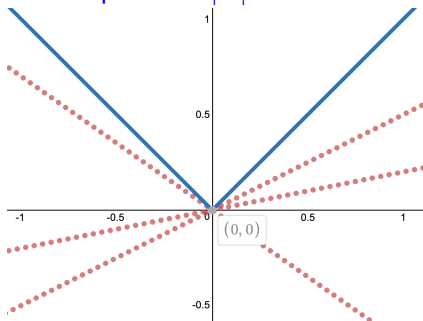
## How do we solve $\nabla f = 0$ when $\nabla f$ doesn't exist?

A **subgradient**  $g \in \mathcal{H}$  of  $f: \mathcal{H} \rightarrow ]-\infty, +\infty]$  at  $x \in \mathcal{H}$  satisfies

$$(\forall y \in \mathcal{H}) \quad \langle y - x \mid g \rangle + f(x) \leq f(y).$$

The **subdifferential**  $\partial f(x) \subset \mathcal{H}$  is the set containing **all subgradients** of  $f$  at  $x$ .

**Example:**  $f = |\cdot|$ : What do we do at zero?  $\partial f(0) = [-1, 1]$



## How do we solve $\nabla f = 0$ when $\nabla f$ doesn't exist?

A **subgradient**  $g \in \mathcal{H}$  of  $f: \mathcal{H} \rightarrow ]-\infty, +\infty]$  at  $x \in \mathcal{H}$  satisfies

$$(\forall y \in \mathcal{H}) \quad \langle y - x \mid g \rangle + f(x) \leq f(y).$$

The **subdifferential**  $\partial f(x) \subset \mathcal{H}$  is the set containing **all subgradients** of  $f$  at  $x$ .

### Fermat's Rule

Let  $x \in \mathcal{H}$ . Then  $x \in \text{Argmin } f$  if and only if  $0 \in \partial f(x)$ .

## How do we solve $\nabla f = 0$ when $\nabla f$ doesn't exist?

A **subgradient**  $g \in \mathcal{H}$  of  $f: \mathcal{H} \rightarrow ]-\infty, +\infty]$  at  $x \in \mathcal{H}$  satisfies

$$(\forall y \in \mathcal{H}) \quad \langle y - x \mid g \rangle + f(x) \leq f(y).$$

The **subdifferential**  $\partial f(x) \subset \mathcal{H}$  is the set containing **all subgradients** of  $f$  at  $x$ .

### Fermat's Rule

Let  $x \in \mathcal{H}$ . Then  $x \in \text{Argmin } f$  if and only if  $0 \in \partial f(x)$ .

*Proof:*

$$\begin{aligned} 0 \in \partial f(x) &\Leftrightarrow (\forall y \in \mathcal{H}) \quad \langle y - x \mid 0 \rangle + f(x) \leq f(y) \\ &\Leftrightarrow (\forall y \in \mathcal{H}) \quad f(x) \leq f(y) \\ &\Leftrightarrow x \in \text{Argmin } f \end{aligned}$$

## How do we solve $\nabla f = 0$ when $\nabla f$ doesn't exist?

A **subgradient**  $g \in \mathcal{H}$  of  $f: \mathcal{H} \rightarrow ]-\infty, +\infty]$  at  $x \in \mathcal{H}$  satisfies

$$(\forall y \in \mathcal{H}) \quad \langle y - x \mid g \rangle + f(x) \leq f(y).$$

The **subdifferential**  $\partial f(x) \subset \mathcal{H}$  is the set containing **all subgradients** of  $f$  at  $x$ .

$\partial f$  is useful for developing **both**  $\underbrace{\text{optimality criterion}}_{0 \in \partial f(x)}$  **and** algorithms.



**Goal: “ $0 \in \partial f(x)$ ”. Which path do we take?**



image: centralldm.es

## Goal: “ $0 \in \partial f(x)$ ”. Which path do we take?

Some *provenly-convergent* (first-order) algorithm classes:

- Subgradient-projections (e.g., in [C. & ZW, *IEEE EUSIPCO*, 2020])
- Proximity operators (e.g., in [C., B., & ZW, *IEEE ICASSP*, 2022], [C. & ZW, *J. Approx. Theory*, 2021], [C. & ZW, *SIAM J. Imaging Sci.*, 2023])
- Conditional Gradient / “Frank-Wolfe” (e.g., in [ZW & P., 2024], [K., P., W., & ZW, *Opt. Methods. Softw.*, 2024])
- Abs-smooth Optimization (e.g., [K., P., W., & ZW, *Opt. Methods. Softw.*, 2024])
- Bundle methods, Barrier methods, Lagrangian methods, ...

## Goal: “ $0 \in \partial f(x)$ ”. Which path do we take?

Some *provenly-convergent* (first-order) algorithm classes:

- Subgradient-projections (e.g., in [C. & ZW, *IEEE EUSIPCO*, 2020])
- **Proximity operators** (e.g., in [C., B., & ZW, *IEEE ICASSP*, 2022], [C. & ZW, *J. Approx. Theory*, 2021], [C. & ZW, *SIAM J. Imaging Sci.*, 2023])
- **Conditional Gradient / “Frank-Wolfe”** (e.g., in [ZW & P., 2024], [K., P., W., & ZW, *Opt. Methods. Softw.*, 2024])
- Abs-smooth Optimization (e.g., [K., P., W., & ZW, *Opt. Methods. Softw.*, 2024])
- Bundle methods, Barrier methods, Lagrangian methods, ...

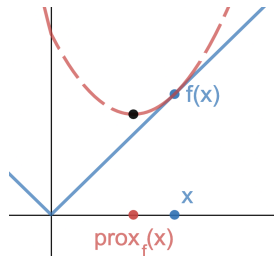
# Theoretically-sound optimization

1. Motivation
2. Background: Theory vs practice
3. Proximity operators: Algorithmic bells and whistles
4. Splitting FW: What if the “usual” tools fail us?
5. More adventuring

# Proximity operators: a new hope

The **proximity operator** of  $f$  at  $x \in \mathcal{H}$  is

$$\text{prox}_f(x) = \underset{u \in \mathcal{H}}{\text{Argmin}} \quad f(u) + \frac{1}{2} \|x - u\|^2$$

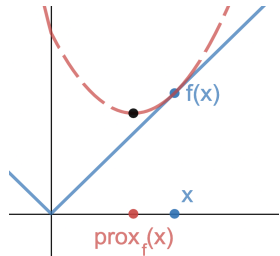


## Proximity operators: a new hope

The **proximity operator** of  $f$  at  $x \in \mathcal{H}$  is

$$\text{prox}_f(x) = \underset{u \in \mathcal{H}}{\text{Argmin}} \quad f(u) + \frac{1}{2} \|x - u\|^2$$

→ For  $f \in \Gamma_0(\mathcal{H})$  and  $x \in \mathcal{H}$ ,  $\text{prox}_f(x)$  is unique.

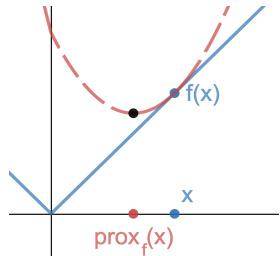


## Proximity operators: a new hope

The **proximity operator** of  $f$  at  $x \in \mathcal{H}$  is

$$\text{prox}_f(x) = \underset{u \in \mathcal{H}}{\text{Argmin}} \quad f(u) + \frac{1}{2} \|x - u\|^2$$

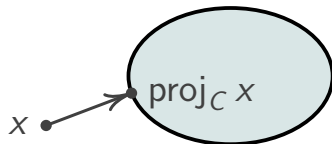
→ For  $f \in \Gamma_0(\mathcal{H})$  and  $x \in \mathcal{H}$ ,  $\text{prox}_f(x)$  is unique. Defines an operator  $\text{prox}_f: \mathcal{H} \rightarrow \mathcal{H}$ .



## Proximity operators: a new hope

The **proximity operator** of  $f$  at  $x \in \mathcal{H}$  is

$$\text{prox}_f(x) = \underset{u \in \mathcal{H}}{\text{Argmin}} \quad f(u) + \frac{1}{2} \|x - u\|^2$$



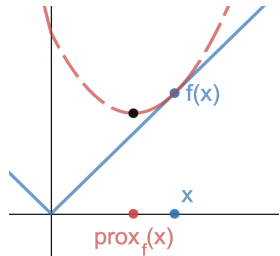
- For  $f \in \Gamma_0(\mathcal{H})$  and  $x \in \mathcal{H}$ ,  $\text{prox}_f(x)$  is unique. Defines an operator  $\text{prox}_f: \mathcal{H} \rightarrow \mathcal{H}$ .
- Projection onto **closed convex set**  $C$ :  $\text{prox}_{\iota_C}(x) = \underset{u \in C}{\text{Argmin}} \|x - u\|^2 = \text{proj}_C x$ .



## Proximity operators: a new hope

The **proximity operator** of  $f$  at  $x \in \mathcal{H}$  is

$$\text{prox}_f(x) = \underset{u \in \mathcal{H}}{\text{Argmin}} \quad f(u) + \frac{1}{2} \|x - u\|^2$$



- For  $f \in \Gamma_0(\mathcal{H})$  and  $x \in \mathcal{H}$ ,  $\text{prox}_f(x)$  is unique. Defines an operator  $\text{prox}_f: \mathcal{H} \rightarrow \mathcal{H}$ .
- Projection onto **closed convex set**  $C$ :  $\text{prox}_{\iota_C}(x) = \underset{u \in C}{\text{Argmin}} \|x - u\|^2 = \text{proj}_C x$ .
- [Martinet, *Fr. Inf. Rech. Oper.*, 1970] (translated / modernized):  
Let  $f \in \Gamma_0(\mathcal{H})$  be such that  $\text{Argmin } f \neq \emptyset$ . Let  $\gamma > 0$ ,  $x_0 \in \mathcal{H}$ , and set

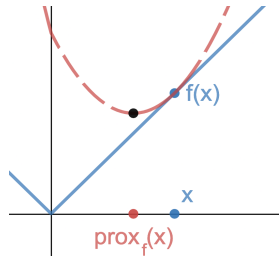
$$x_{n+1} = \text{prox}_{\gamma f} x_n.$$

Then  $(x_n)_{n \in \mathbb{N}} \rightharpoonup x^* \in \text{Argmin } f$ .

## Proximity operators: a new hope

The **proximity operator** of  $f$  at  $x \in \mathcal{H}$  is

$$\text{prox}_f(x) = \underset{u \in \mathcal{H}}{\text{Argmin}} \quad f(u) + \frac{1}{2} \|x - u\|^2$$



→ For  $f \in \Gamma_0(\mathcal{H})$  and  $x \in \mathcal{H}$ ,  $\text{prox}_f(x)$  is unique. Defines an operator  $\text{prox}_f: \mathcal{H} \rightarrow \mathcal{H}$ .

→ Projection onto **closed convex set**  $C$ :  $\text{prox}_{\iota_C}(x) = \underset{u \in C}{\text{Argmin}} \|x - u\|^2 = \text{proj}_C x$ .

→ [Martinet, *Fr. Inf. Rech. Oper.*, 1970] (translated / modernized):

Let  $f \in \Gamma_0(\mathcal{H})$  be such that  $\text{Argmin } f \neq \emptyset$ . Let  $\gamma > 0$ ,  $x_0 \in \mathcal{H}$ , and set

$$x_{n+1} = \text{prox}_{\gamma f} x_n.$$

Then  $(x_n)_{n \in \mathbb{N}} \rightharpoonup x^* \in \text{Argmin } f$ . **Issue:**  $\text{prox}_f$  might be hard to compute.

# Evolution of prox-based algorithms

$$f = \sum_{1 \leq i \leq m} f_i$$

$(\text{prox}_{f_i})_{1 \leq i \leq m}$  are simpler than  $\text{prox}_f$ .

# Evolution of prox-based algorithms

$$f = \sum_{1 \leq i \leq m} f_i$$

$(\text{prox}_{f_i})_{1 \leq i \leq m}$  are simpler than  $\text{prox}_f$ .

Open-source repo's:

Python/Matlab:

[proximity-operator.net](https://proximity-operator.net),

Julia:

[ProximalOperators.jl](https://github.com/ProximalOperators/ProximalOperators.jl) (Github)

# Evolution of prox-based algorithms

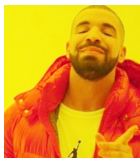
$$f = \sum_{1 \leq i \leq m} f_i$$

$(\text{prox}_{f_i})_{1 \leq i \leq m}$  are simpler than  $\text{prox}_f$ .

Algorithmic “Splitting” mentality:



$\text{prox}_f$



$\text{prox}_{f_1}, \dots, \text{prox}_{f_m}$

# Evolution of prox-based algorithms

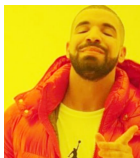
$$f = \sum_{1 \leq i \leq m} f_i$$

$(\text{prox}_{f_i})_{1 \leq i \leq m}$  are simpler than  $\text{prox}_f$ .

Algorithmic “Splitting” mentality:



$\text{prox}_f$



$\text{prox}_{f_1}, \dots, \text{prox}_{f_m}$

## Pseudocode for many proximal splitting algorithms\*

**Require:** Point  $x_0 \in \mathcal{H}$ , objective function  $f$

```
1: for  $n = 0, 1$  to ... do
2:   # Preprocess  $x_n$ 
3:   for  $i = 1$  to  $m$  do
4:      $y_{i,n+1} \leftarrow$  evaluation of  $\text{prox}_{f_i}(\cdot)$ 
5:   end for
6:    $x_{n+1} \leftarrow$  combine  $(y_{i,n+1})_{1 \leq i \leq m}$ 
7: end for
```

\*:e.g., Douglas-Rachford, ADMM, Chambolle-Pock, Forward-backward, Augmented Lagrangian, ...

# Evolution of prox-based algorithms

$$f = \sum_{1 \leq i \leq m} f_i$$

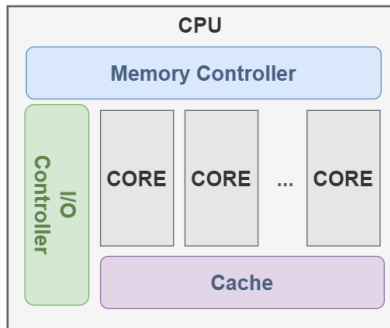


image: baeldung.com

## Pseudocode for many proximal splitting algorithms\*

**Require:** Point  $x_0 \in \mathcal{H}$ , objective function  $f$

- 1: **for**  $n = 0, 1$  **to**  $\dots$  **do**
- 2:   # Preprocess  $x_n$
- 3:   **for**  $i = 1$  **to**  $m$  **do**
- 4:      $y_{i,n+1} \leftarrow$  evaluation of  $\text{prox}_{f_i}(\cdot)$
- 5:   **end for**
- 6:    $x_{n+1} \leftarrow$  combine  $(y_{i,n+1})_{1 \leq i \leq m}$
- 7: **end for**

\*:e.g., Douglas-Rachford, ADMM, Chambolle-Pock, Forward-backward, Augmented Lagrangian, ...

# Evolution of prox-based algorithms

$$f = \sum_{1 \leq i \leq m} f_i$$

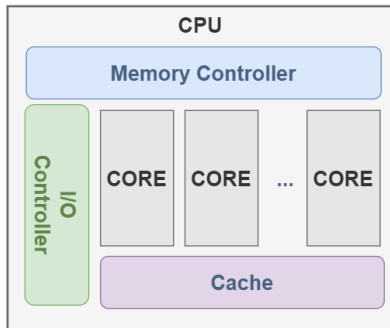


image: baeldung.com

## Pseudocode for many proximal splitting algorithms\*

**Require:** Point  $x_0 \in \mathcal{H}$ , objective function  $f$

- 1: **for**  $n = 0, 1$  **to**  $\dots$  **do**
- 2:   # Preprocess  $x_n$
- 3:   **for**  $i = 1$  **to**  $m$  **do**
- 4:      $y_{i,n+1} \leftarrow$  **evaluation** of  $\text{prox}_{f_i}(\cdot)$
- 5:   **end for**
- 6:    $x_{n+1} \leftarrow$  **combine**  $(y_{i,n+1})_{1 \leq i \leq m}$
- 7: **end for**

\*:e.g., Douglas-Rachford, ADMM, Chambolle-Pock, Forward-backward, Augmented Lagrangian, ...



# Evolution of prox-based algorithms

$$f = \sum_{1 \leq i \leq m} f_i$$

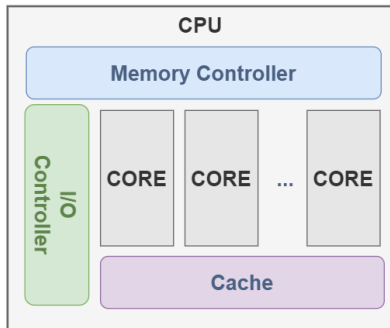


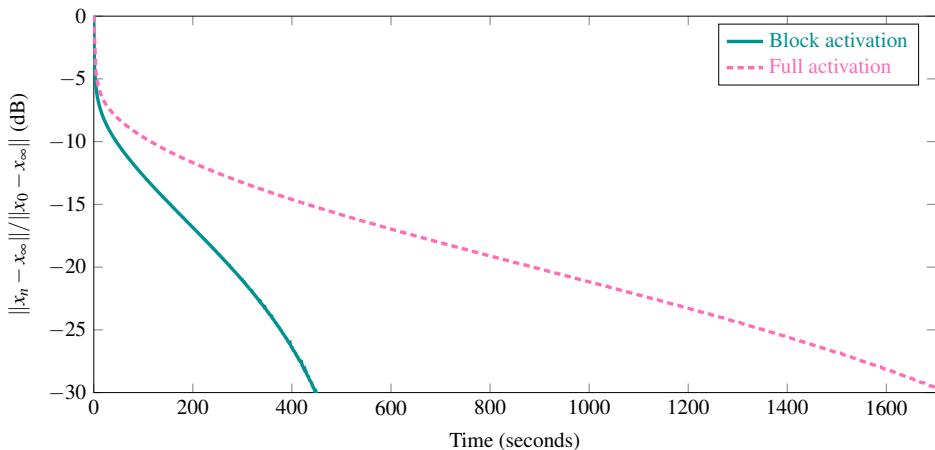
image: baeldung.com

## Pseudocode for block-iterative proximal splitting algorithms

**Require:** Point  $x_0 \in \mathcal{H}$ , objective function  $f$

```
1: for  $n = 0, 1$  to ... do
2:   # Preprocess  $x_n$ ; Select  $I_n \subset \{1, \dots, m\}$ 
3:   for  $i = 1$  to  $m$  do
4:     if  $i \in I_n$  then
5:        $y_{i,n+1} \leftarrow$  evaluation of  $\text{prox}_{f_i}(\cdot)$ 
6:     else
7:        $y_{i,n+1} \leftarrow y_{i,n}$ 
8:     end if
9:   end for
10:   $x_{n+1} \leftarrow$  combine  $(y_{i,n+1})_{1 \leq i \leq m}$ 
11: end for
```

## Block activation for image recovery ( $m = 2$ )



[C. & ZW, *SIAM J. Imaging Sci.*, 2022] Relative error for **full-activation** ( $l_n = I$ ) versus **block activation**:

$$l_n = \begin{cases} \{1, 2\}, & \text{if } n \equiv 0 \pmod{5}; \\ \{2\}, & \text{if } n \not\equiv 0 \pmod{5}. \end{cases}$$

# Evolution of prox-based algorithms

$$f = \sum_{1 \leq i \leq m} f_i$$

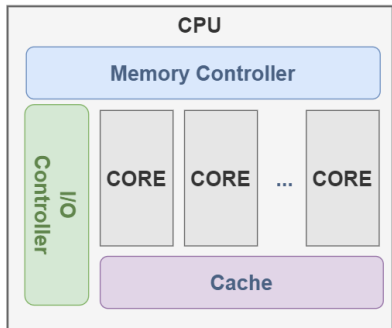


image: baeldung.com

## Pseudocode for block-iterative proximal splitting algorithms\*

**Require:** Point  $x_0 \in \mathcal{H}$ , objective function  $f$

```
1: for  $n = 0, 1$  to ... do
2:   # Preprocess  $x_n$ ; Select  $I_n \subset \{1, \dots, m\}$ 
3:   for  $i = 1$  to  $m$  do
4:     if  $i \in I_n$  then
5:        $y_{i,n+1} \leftarrow$  evaluation of  $\text{prox}_{f_i}(\cdot)$ 
6:     else
7:        $y_{i,n+1} \leftarrow y_{i,n}$ 
8:     end if
9:   end for
10:   $x_{n+1} \leftarrow$  combine  $(y_{i,n+1})_{1 \leq i \leq m}$ 
11: end for
```

# Evolution of prox-based algorithms

$$f = \sum_{1 \leq i \leq m} f_i$$

Two (currently separate) approaches:

→ Asynchronous updates

→ Extrapolated updates

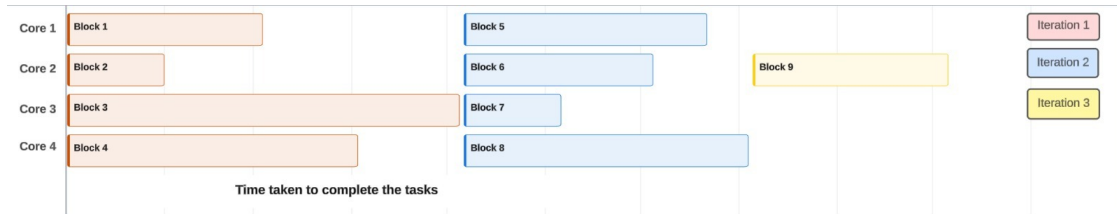
## Pseudocode for block-iterative proximal splitting algorithms\*

**Require:** Point  $x_0 \in \mathcal{H}$ , objective function  $f$

```
1: for  $n = 0, 1$  to ... do
2:   # Preprocess  $x_n$ ; Select  $I_n \subset \{1, \dots, m\}$ 
3:   for  $i = 1$  to  $m$  do
4:     if  $i \in I_n$  then
5:        $y_{i,n+1} \leftarrow$  evaluation of  $\text{prox}_{f_i}(\cdot)$ 
6:     else
7:        $y_{i,n+1} \leftarrow y_{i,n}$ 
8:     end if
9:   end for
10:   $x_{n+1} \leftarrow$  combine  $(y_{i,n+1})_{1 \leq i \leq m}$ 
11: end for
```

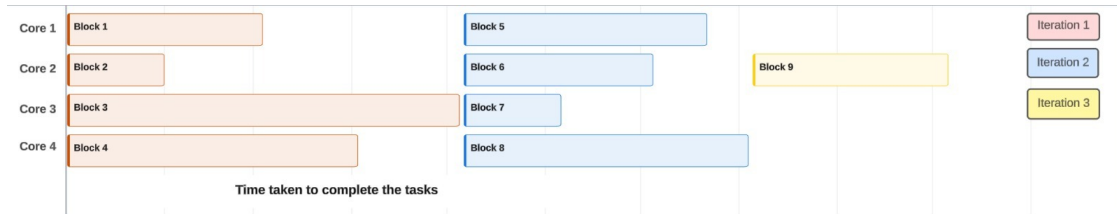
# Evolution of prox-based algorithms

## Parallel and Synchronous

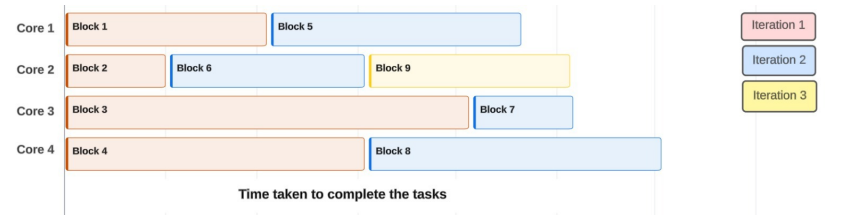


# Evolution of prox-based algorithms

## Parallel and Synchronous



## Parallel and Asynchronous



# Asynchrony: Projective Splitting Algorithms

[Eckstein & Svaiter, *Math Prog. A*, 2008]:

Coined “Projective splitting” (synchronous, not block-iterative).

[Combettes & Eckstein, *Math Prog. B*, 2018]:

Convergence proof with **asynchronous block-iterative** updates!

# Asynchrony: Projective Splitting Algorithms

[Eckstein & Svaiter, *Math Prog. A*, 2008]:

Coined “Projective splitting” (synchronous, not block-iterative).

[Combettes & Eckstein, *Math Prog. B*, 2018]:

Convergence proof with **asynchronous block-iterative** updates!

[Combettes, Búi, & ZW, *IEEE ICASSP*, 2022]:

Numerical analysis (space and time complexity). For ML (training classifiers) and image processing, works better than other algorithms in its class.



# Asynchrony: Projective Splitting Algorithms

[Eckstein & Svaiter, *Math Prog. A*, 2008]:

Coined “Projective splitting” (synchronous, not block-iterative).

[Combettes & Eckstein, *Math Prog. B*, 2018]:

Convergence proof with **asynchronous block-iterative** updates!

[Combettes, Búi, & ZW, *IEEE ICASSP*, 2022]:

Numerical analysis (space and time complexity). For ML (training classifiers) and image processing, works better than other algorithms in its class.

Only studied **synchronous** case!

# Asynchrony: Projective Splitting Algorithms

[Eckstein & Svaiter, *Math Prog. A*, 2008]:

Coined “Projective splitting” (synchronous, not block-iterative).

[Combettes & Eckstein, *Math Prog. B*, 2018]:

Convergence proof with **asynchronous block-iterative** updates!

[Combettes, Búi, & ZW, *IEEE ICASSP*, 2022]:

Numerical analysis (space and time complexity). For ML (training classifiers) and image processing, works better than other algorithms in its class.

Only studied **synchronous** case!

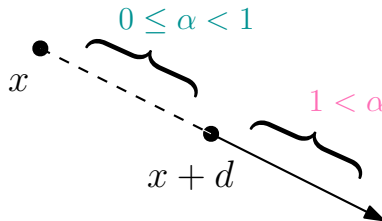
Dua, Goel, Sharma, & ZW (ongoing work):

Analysis and experimentation for **asynchronous** case. AsyncProx.jl in development.

## Over-relaxation (extrapolation)

Given a “descent direction”  $d$  (e.g., combination of  $(y_{i,n+1})_{1 \leq i \leq m}$ ) from  $x \in \mathcal{H}$ ,

$$x_+ = (1 - \alpha)x + \alpha(x + d)$$



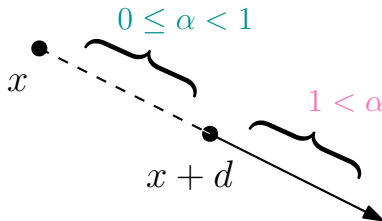
## Over-relaxation (extrapolation)

Given a “descent direction”  $d$  (e.g., combination of  $(y_{i,n+1})_{1 \leq i \leq m}$ ) from  $x \in \mathcal{H}$ ,

$$x_+ = (1 - \alpha)x + \alpha(x + d)$$

$0 \leq \alpha < 1$ : under-relaxation

$1 < \alpha$ : over-relaxation



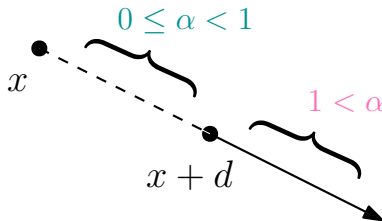
## Over-relaxation (extrapolation)

Given a “descent direction”  $d$  (e.g., combination of  $(y_{i,n+1})_{1 \leq i \leq m}$ ) from  $x \in \mathcal{H}$ ,

$$x_+ = (1 - \alpha)x + \alpha(x + d)$$

$0 \leq \alpha < 1$ : under-relaxation

$1 < \alpha$ : over-relaxation



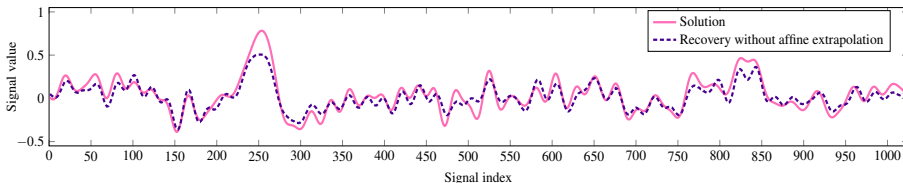
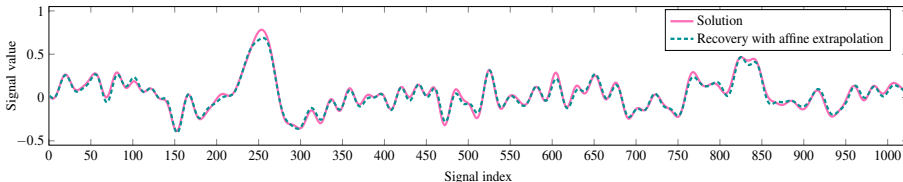
$$\begin{aligned} x_{n+1} &= x_n - \alpha_n \frac{1}{L} \nabla f(x_n), & \text{where } \alpha_n &\in [\varepsilon, 2 - \varepsilon] & \quad (\text{GD}) \\ &= (1 - \alpha_n)x_n + \alpha_n \left( x_n - \frac{1}{L} \nabla f(x_n) \right) \end{aligned}$$

# Over-relaxation for fixed-point problems

[Combettes & ZW, *J. Approx. Theory*, 2021]:

A strongly-convergent algorithm with [affine-convex extrapolation](#).

**Example:** EEG data (minimal-norm solution to an ill-posed nonlinear inverse problem;  
recovery after 1000 iterations, < 1 minute.)



# Theoretically-sound optimization

1. Motivation
2. Background: Theory vs practice
3. Proximity operators: Algorithmic bells and whistles
4. Splitting FW: What if the “usual” tools fail us?
5. More adventuring

# They stole my horse!

## Splitting problem setup

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

$$\text{minimize } f(x) \text{ subject to } x \in \bigcap_{1 \leq i \leq m} C_i, \quad (\star)$$

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



# They stole my horse!

## Splitting problem setup

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

$$\text{minimize } f(x) \text{ subject to } x \in \bigcap_{1 \leq i \leq m} C_i, \quad (\star)$$

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

What if  $(\text{proj}_{C_i})_{1 \leq i \leq m}$  are too expensive?  $\leftarrow$  e.g., in high-dimensional settings!

# A spark of inspiration

Frank-Wolfe / “Conditional gradient” alg. [Naval Res. Logist. Quart., 1956]

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

minimize  $f(x)$  subject to  $x \in C$

Instead of projecting, use a **linear minimization oracle** of  $C$ ,

$$\text{LMO}_C: y \mapsto p \in \text{Argmin}_{x \in C} \langle y | x \rangle \quad (\text{LMO})$$

$$x_{n+1} = x_n + \frac{1}{n+1} \left( \text{LMO}_C(\nabla f(x_n)) - x_n \right)$$



Marguerite Frank



Philip Wolfe

# efficiency(LMOs) == efficiency(lembus bread)

[Combettes/Pokutta, '21]: For many constraints,  $C$ ,  $\text{proj}_C$  is **more expensive** than  $\text{LMO}_C$ .  
(e.g., nuclear norm ball,  $\ell_1$  ball, probability simplex, Birkhoff polytope, general LP, ...)

# efficiency(LMOs) == efficiency(lembus bread)

[Combettes/Pokutta, '21]: For many constraints,  $C$ ,  $\text{proj}_C$  is **more expensive** than  $\text{LMO}_C$ .  
(e.g., nuclear norm ball,  $\ell_1$  ball, probability simplex, Birkhoff polytope, general LP, ...)

## Example: Nuclear norm ball

For  $x \in \mathbb{R}^{n \times n}$

$$\|x\|_{nuc} = \sum_{1 \leq i \leq n} \sigma_i(x).$$

For  $\beta \geq 0$ ,  $C = \{x \in \mathbb{R}^{n \times n} \mid \|x\|_{nuc} \leq \beta\}$ ,

# efficiency(LMOs) == efficiency(lembus bread)

[Combettes/Pokutta, '21]: For many constraints,  $C$ ,  $\text{proj}_C$  is **more expensive** than  $\text{LMO}_C$ .  
(e.g., nuclear norm ball,  $\ell_1$  ball, probability simplex, Birkhoff polytope, general LP, ...)

## Example: Nuclear norm ball

For  $x \in \mathbb{R}^{n \times n}$

$$\|x\|_{nuc} = \sum_{1 \leq i \leq n} \sigma_i(x).$$

For  $\beta \geq 0$ ,  $C = \{x \in \mathbb{R}^{n \times n} \mid \|x\|_{nuc} \leq \beta\}$ ,

$\text{proj}_C(x)$ : requires a **full SVD**!

$(\sigma_1, \dots, \sigma_n, U, V)$ , where  $x = U \text{diag}(\sigma_1, \dots, \sigma_n) V^\top$

Full SVD

$n = 500$ : 0.11 sec

$n = 1000$ : 0.47 sec

$n = 2000$ : 4.87 sec

# efficiency(LMOs) == efficiency(lembus bread)

[Combettes/Pokutta, '21]: For many constraints,  $C$ ,  $\text{proj}_C$  is **more expensive** than  $\text{LMO}_C$ .  
(e.g., nuclear norm ball,  $\ell_1$  ball, probability simplex, Birkhoff polytope, general LP, ...)

## Example: Nuclear norm ball

For  $x \in \mathbb{R}^{n \times n}$

$$\|x\|_{nuc} = \sum_{1 \leq i \leq n} \sigma_i(x).$$

For  $\beta \geq 0$ ,  $C = \{x \in \mathbb{R}^{n \times n} \mid \|x\|_{nuc} \leq \beta\}$ ,

$\text{proj}_C(x)$ : requires a **full SVD**!

$(\sigma_1, \dots, \sigma_n, U, V)$ , where  $x = U \text{diag}(\sigma_1, \dots, \sigma_n) V^\top$

$\text{LMO}_C(x)$ : requires only **first singular value/vectors**  
 $(\sigma_1, U_1, V_1^\top)$

### Full SVD

$n = 500$ : 0.11 sec

$n = 1000$ : 0.47 sec

$n = 2000$ : 4.87 sec

### Just $(\sigma_1, U_1, V_1^\top)$

$n = 500$ : 0.0081 sec

$n = 1000$ : 0.056 sec

$n = 2000$ : 0.638 sec

## Sure would be nice if we could use splitting...

$$\text{minimize } f(x) \text{ subject to } x \in \bigcap_{1 \leq i \leq m} C_i, \quad (\star)$$

## Sure would be nice if we could use splitting...

$$\text{minimize } f(x) \text{ subject to } x \in \bigcap_{1 \leq i \leq m} C_i, \quad (\star)$$

**Issue:** Computing the LMO for  $\bigcap_{1 \leq i \leq m} C_i$  is prohibitively expensive.



## Sure would be nice if we could use splitting...

$$\text{minimize } f(x) \text{ subject to } x \in \bigcap_{1 \leq i \leq m} C_i, \quad (\star)$$

**Issue:** Computing the LMO for  $\bigcap_{1 \leq i \leq m} C_i$  is prohibitively expensive.

$(\text{LMO}_{C_i})_{1 \leq i \leq m}$  are easier to evaluate (e.g., repository: FrankWolfe.jl)

## Sure would be nice if we could use splitting...

$$\text{minimize } f(x) \text{ subject to } x \in \bigcap_{1 \leq i \leq m} C_i, \quad (\star)$$

**Issue:** Computing the LMO for  $\bigcap_{1 \leq i \leq m} C_i$  is prohibitively expensive.

$(\text{LMO}_{C_i})_{1 \leq i \leq m}$  are easier to evaluate (e.g., repository: FrankWolfe.jl)

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

Use  $\text{LMO}_{C_1}, \text{LMO}_{C_2}, \dots$  instead of  $\text{LMO}_{(\bigcap_{1 \leq i \leq m} C_i)}$

## Sure would be nice if we could use splitting...

$$\text{minimize } f(x) \text{ subject to } x \in \bigcap_{1 \leq i \leq m} C_i, \quad (\star)$$

**Issue:** Computing the LMO for  $\bigcap_{1 \leq i \leq m} C_i$  is prohibitively expensive.

$(\text{LMO}_{C_i})_{1 \leq i \leq m}$  are easier to evaluate (e.g., repository: FrankWolfe.jl)

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

Use  $\text{LMO}_{C_1}, \text{LMO}_{C_2}, \dots$  instead of  $\text{LMO}_{(\bigcap_{1 \leq i \leq m} C_i)}$

Relatively little has been done in this field.

## Sure would be nice if we could use splitting...

$$\text{minimize } f(x) \text{ subject to } x \in \bigcap_{1 \leq i \leq m} C_i, \quad (\star)$$

**Issue:** Computing the LMO for  $\bigcap_{1 \leq i \leq m} C_i$  is prohibitively expensive.

$(\text{LMO}_{C_i})_{1 \leq i \leq m}$  are easier to evaluate (e.g., repository: FrankWolfe.jl)

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

Use  $\text{LMO}_{C_1}, \text{LMO}_{C_2}, \dots$  instead of  $\text{LMO}_{(\bigcap_{1 \leq i \leq m} C_i)}$

Relatively little has been done in this field.

→ Unlike projections, LMOs are discontinuous.

## Sure would be nice if we could use splitting...

$$\text{minimize } f(x) \text{ subject to } x \in \bigcap_{1 \leq i \leq m} C_i, \quad (\star)$$

**Issue:** Computing the LMO for  $\bigcap_{1 \leq i \leq m} C_i$  is prohibitively expensive.

$(\text{LMO}_{C_i})_{1 \leq i \leq m}$  are easier to evaluate (e.g., repository: FrankWolfe.jl)

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

Use  $\text{LMO}_{C_1}, \text{LMO}_{C_2}, \dots$  instead of  $\text{LMO}_{(\bigcap_{1 \leq i \leq m} C_i)}$

Relatively little has been done in this field.

→ Unlike projections, LMOs are discontinuous.

→ “State-of-the-art” relies on inexact prox-based algorithms.

# It's aliive!

## [ZW & P, 2024]: Split Conditional Gradient Algorithm

**Require:** Point  $x_0 \in \frac{1}{m} \sum_{1 \leq i \leq m} C_i$ , smooth function  $f$

```
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:      $v_t^i \leftarrow \text{LMO}_{C_i}(g_t + \lambda_t(x_t^i - x_t))$   
7:      $x_{t+1}^i \leftarrow x_t^i + \gamma_t(v_t^i - x_t^i)$   
8:   end for  
9:    $x_{t+1} \leftarrow \frac{1}{m} \sum_{1 \leq i \leq m} x_{t+1}^i$   
10: end for
```

Practical advantages:

- Uses individual LMOs
- Lowest-known # LMO calls:  
one LMO per set (per iteration)

# It's aliive!

## [ZW & P, 2024]: Split Conditional Gradient Algorithm

**Require:** Point  $x_0 \in \frac{1}{m} \sum_{1 \leq i \leq m} C_i$ , smooth function  $f$

```
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:      $v_t^i \leftarrow \text{LMO}_{C_i}(g_t + \lambda_t(x_t^i - x_t))$ 
7:      $x_{t+1}^i \leftarrow x_t^i + \gamma_t(v_t^i - x_t^i)$ 
8:   end for
9:    $x_{t+1} \leftarrow \frac{1}{m} \sum_{1 \leq i \leq m} x_{t+1}^i$ 
10: end for
```

Practical advantages:

- Uses individual LMOs
- Lowest-known # LMO calls:  
one LMO per set (per iteration)

**Q:** Does it actually solve  $(\star)$ ?

**A:** Yes.

$\gamma_t = \mathcal{O}(1/\sqrt{t})$  and  
 $\lambda_t = \mathcal{O}(\ln t)$  work.  
(whether or not  $f$  is convex).

# Convergence

## Theorem ([ZW & P., 2024])

Let  $f$  be  $L$ -smooth and let  $(C_i)_{1 \leq i \leq m}$  be nonempty compact convex subsets of  $\mathcal{H}$  such that  $\bigcap_{1 \leq i \leq m} C_i \neq \emptyset$ . Let  $\lambda_0 > 0$  and  $\lambda_{t+1} = \lambda_t + (\sqrt{t} + 2)^{-2}$  and  $\gamma_t = 2/(\sqrt{t} + 2)$ . If  $f$  is **convex**, then

$$f\left(\frac{1}{m} \sum_{1 \leq i \leq m} \mathbf{x}_t^i\right) \rightarrow \inf_{x \in \bigcap_{1 \leq i \leq m} C_i} f(x)$$

and every accumulation point  $\mathbf{x}_\infty$  of  $(\mathbf{x}_t)_{t \in \mathbb{N}}$  produces a solution

$$\frac{1}{m} \sum_{1 \leq i \leq m} \mathbf{x}_\infty^i \in \bigcap_{1 \leq i \leq m} C_i \text{ such that } f\left(\frac{1}{m} \sum_{1 \leq i \leq m} \mathbf{x}_\infty^i\right) = \inf_{x \in \bigcap_{1 \leq i \leq m} C_i} f(x).$$

**Nonconvex** convergence results too: arXiv:2311.05381



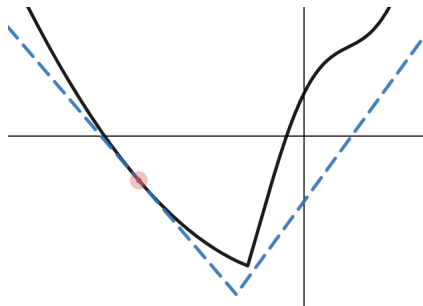
# Theoretically-sound optimization

1. Motivation
2. Background: Theory vs practice
3. Proximity operators: Algorithmic bells and whistles
4. Splitting FW: What if the “usual” tools fail us?
- 5. More adventuring**

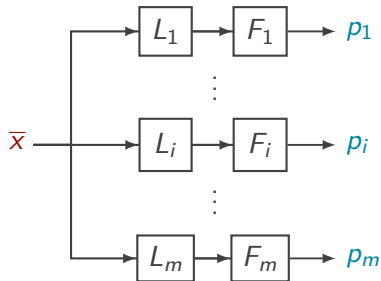
# Abs-smooth optimization

**Abs-smooth** functions  $f$  include compositions of smooth functions, max, min, and  $|\cdot|$ . ...

- Loss functions for multilayer **Neural Networks** with smooth and/or ReLU activation.
- Allows one to find a local minimizer on **non-convex** objective functions!
- Future work: Improve scalability.



## Nonlinear inverse problems

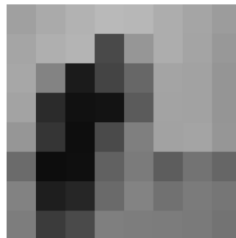


Given  $(p_i)_{1 \leq i \leq m}$ , find  $x^* \in \mathcal{H}$  such that

$$(\forall i \in \{1, \dots, m\}) \quad F_i(L_i x^*) = p_i$$

$L_i$  are bounded linear operators, and  $F_i \approx$  proximity operators.

→ To-do: Stability analysis



# Outlook and future work

- Improved convergence rates and acceleration
- **Block-iterative** Frank-Wolfe algorithms
- Efficient prox/ LMO algorithms
- ...



## Outlook and future work

- Improved convergence rates and acceleration
- **Block-iterative** Frank-Wolfe algorithms
- Efficient prox/ LMO algorithms
- ...

Potential collaborators: Hala Nelson, Minah Oh, Roger Thelwell, and more!

For students:






Proofs (MATH 245), sequences and series (236), gradients (237), linear algebra (238/300/434), optimization theory (340), coding experience (248/250/448/449), analysis (410/411).

REUs & Grad School in Berlin: [iol.zib.de](http://iol.zib.de) (Opt/ML), [math-berlin.de](http://math-berlin.de)



**Thank you for your attention!**

# References

-  M. N. Búi, P. L. Combettes and ZW, [Block-activated algorithms for multicomponent fully nonsmooth minimization](#)  
*Proc. IEEE Int. Conf. Acoust. Speech Signal Process.*, pp 5428–5432, 2022.
-  P. L. Combettes, A. M. McDonald, C. A. Micchelli, and M. Pontil, [Learning with optimal interpolation norms](#)  
*Numer. Algorithms*, vol. 81, no. 2, pp. 695–717, 2019.
-  P. L. Combettes and ZW, [Signal recovery from inconsistent nonlinear observations](#)  
*Proc. IEEE Int. Conf. Acoust. Speech Signal Process.*, pp 5872–5876, 2022.
-  P. L. Combettes and ZW, [A fixed point framework for recovering signals from nonlinear transformations](#)  
*Proc. IEEE Eur. Signal Proc.*, pp 2120–2124, 2020.
-  P. L. Combettes and ZW, [Reconstruction of functions from prescribed proximal points](#)  
*J. Approx. Theory*, vol. 268, no. 105606, 2021

# References



P. L. Combettes and ZW, [A variational inequality model for the construction of signals from inconsistent nonlinear equations](#)

*SIAM J. Imaging Sci.*, vol. 15, no. 1, pp. 84–109, 2022



C. Combettes and S. Pokutta, Complexity of linear minimization and projection on some sets  
*Oper. Res. Lett.*, vol. 49, no. 4, pp. 565–571, 2021



M. Frank and P. Wolfe, An algorithm for quadratic programming  
*Naval Res. Logist. Quart.*, vol. 3, iss. 1–2, pp. 95–110, 1956



T. Kreimeier, S. Pokutta, A. Walther, and ZW, [On a Frank-Wolfe approach for abs-smooth functions](#)




*Opt. Methods and Softw.*, DOI: 10.1080/10556788.2023.2296985, 2024



B. Martinet, régularisation d'inéquations variationnelles par approximations successives  
*Fr. Inf. Rech. Oper.*, Série rouge, 4 (R3):154–158, 1970.



# References

-  C. S. Sartori and L. S. Buriol, A study on the pickup and delivery problem with time windows: Matheuristics and new instances  
*Comput. Oper. Res.*, vol. 124, art. 105065, 2020.
-  N. Torelli, D. Papp, and J. Unkelbach, Spatiotemporal fractionation schemes for stereotactic radiosurgery of multiple brain metastases  
*Med. Phys.*, vol. 50, no. 8, pp. 5095-5114, 2023.
-  ZW and S. Pokutta, [Splitting the conditional gradient algorithm](#)  
*arXiv:2311.05381*, 2024

# Supplement

$$f = \sum_{1 \leq i \leq m} f_i$$

$(\text{prox}_{f_i})_{1 \leq i \leq m}$  are simpler than  $\text{prox}_f$ .

e.g., Find  $x \in \mathbb{S}_+^n \cap [\alpha, \beta]^{N \times N}$

minimize  $\iota_{\mathbb{S}_+^n}(x) + \iota_{[\alpha, \beta]^N}(x)$

$\text{prox}_f$  is intractable.

$\left. \begin{array}{l} \text{prox}_{f_1} = \text{proj}_{[\alpha, \beta]^N} \\ \text{prox}_{f_2} = \text{proj}_{\mathbb{S}_+^n} \end{array} \right\}$  known in closed-form.