

Proximal Algorithms for Large-scale Convex Nonsmooth Optimization

Laurent Condat



Visual Computing Center & AI Initiative
King Abdullah Univ. of Science and Technology
(KAUST), Saudi Arabia

Optimization

“Nothing takes place in the world whose meaning is not that of some maximum or minimum”

— **Leonhard Euler** ~1750

Optimization

Find $x^* \in \arg \min_{x \in \mathcal{X}} \psi(x)$

Optimization

Find $x^* \in \arg \min_{x \in \mathcal{X}} \Psi(x)$

\mathcal{X} is a real Hilbert space

Optimization

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \psi(x) = \sum_{m=1}^M g_m(L_m x)$$

with

- linear operators $L_m : \mathcal{X} \rightarrow \mathcal{U}_m$
- real Hilbert spaces $\mathcal{X}, \mathcal{U}_m$
- functions g_m

Optimization

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \psi(x) = \sum_{m=1}^M g_m(L_m x)$$

$$\text{Example: } \psi(x) = \|Ax - y\|^2 + \|x\|_1$$

Optimization

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \psi(x) = \sum_{m=1}^M g_m(L_m x)$$

$$\text{Example: } \psi(x) = \|Ax - y\|^2 + R(x)$$

Motivation: image processing

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \Psi(x) = \sum_{m=1}^M g_m(L_m x)$$

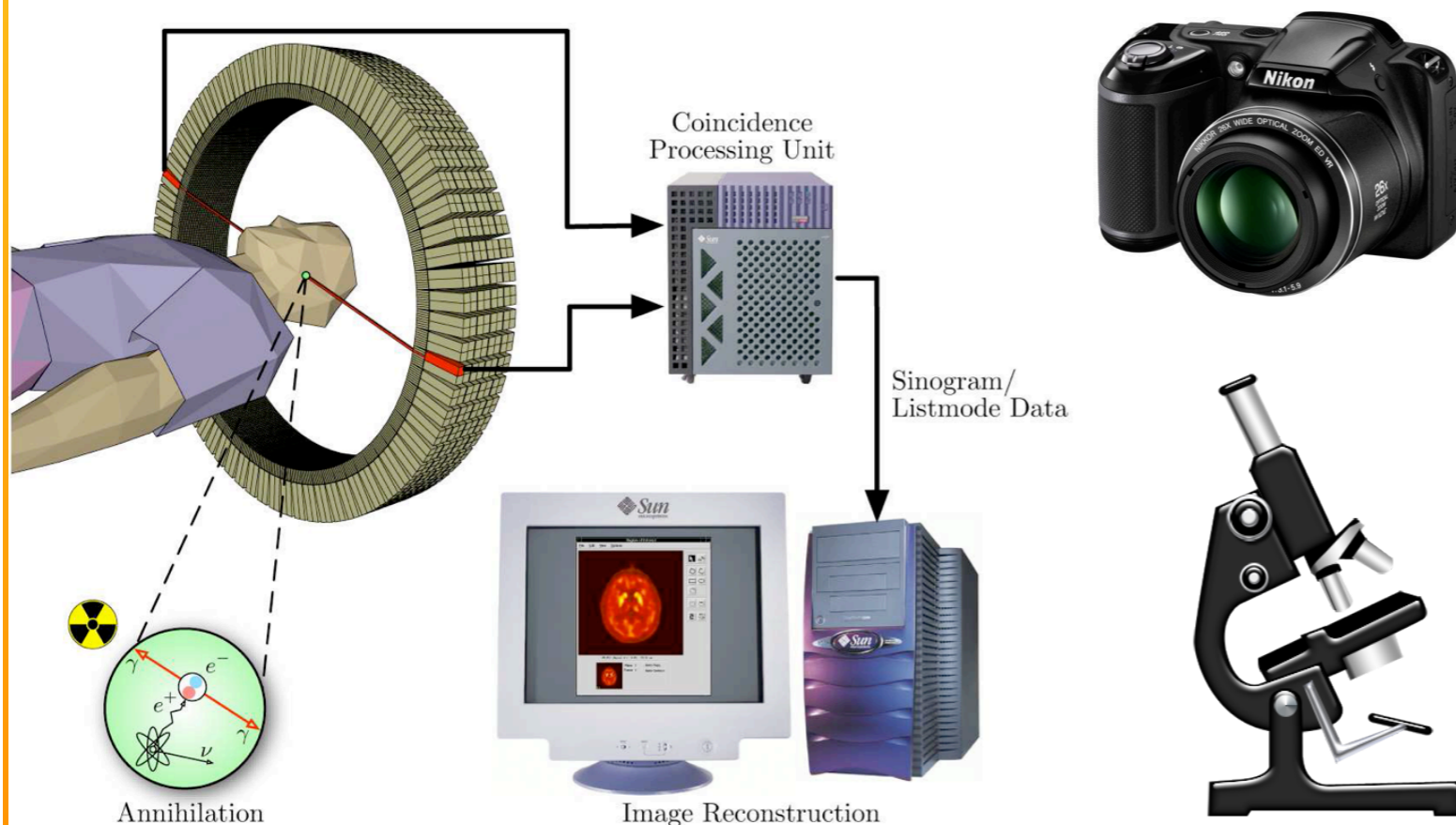
$$\text{Example: } \Psi(x) = \|Ax - y\|^2 + R(x)$$



Motivation: image processing

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \Psi(x) = \sum_{m=1}^M g_m(L_m x)$$

$$\text{Example: } \Psi(x) = \|Ax - y\|^2 + R(x)$$



LC, "Discrete total variation: New definition and minimization", 2017

LC, "A generic proximal algorithm...", 2014

Convex optimization

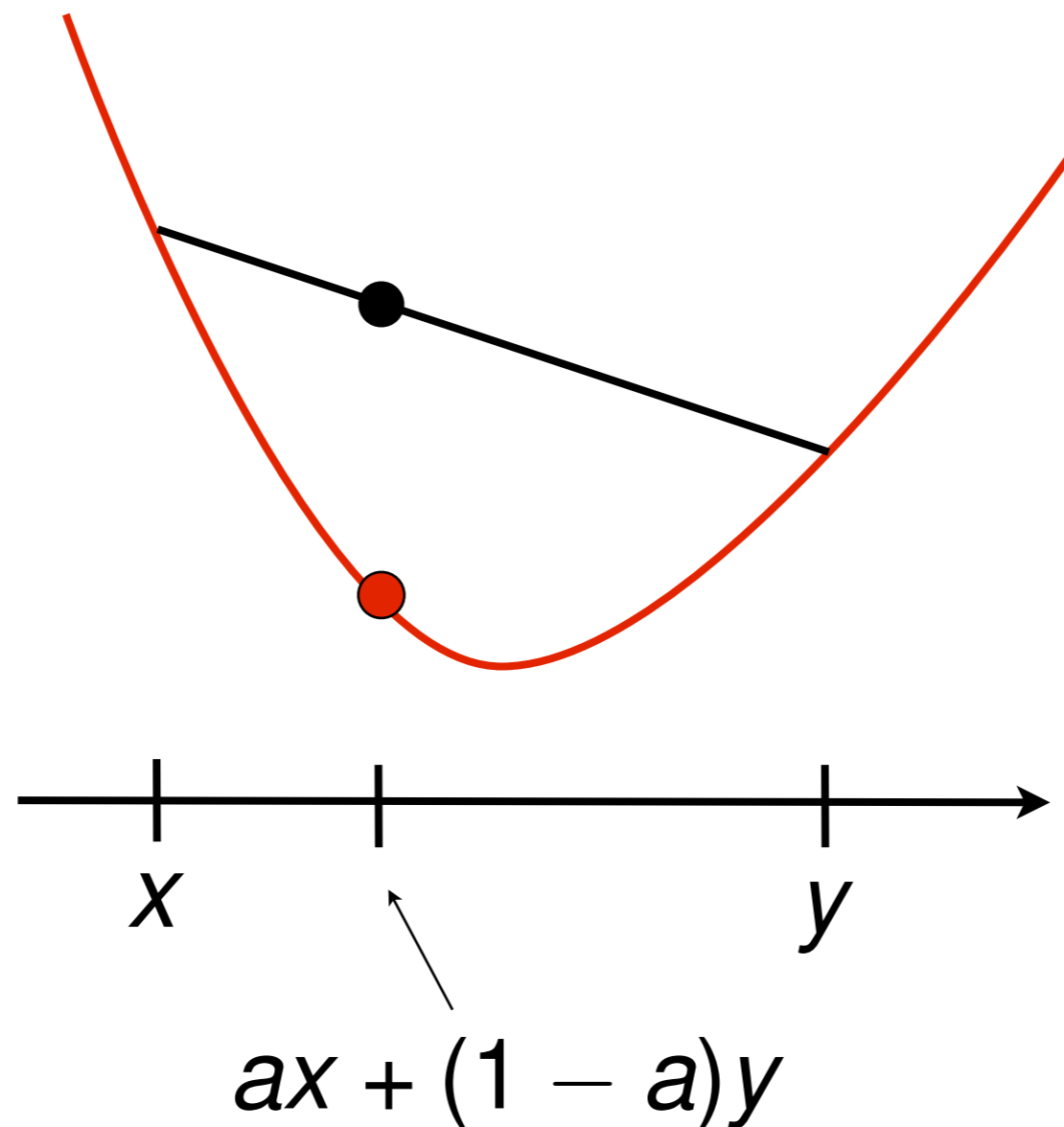
$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \psi(x) = \sum_{m=1}^M g_m(L_m x)$$

with

- linear operators $L_m : \mathcal{X} \rightarrow \mathcal{U}_m$
- real Hilbert spaces $\mathcal{X}, \mathcal{U}_m$
- **convex** functions $g_m : \mathcal{U}_m \rightarrow \mathbb{R} \cup \{+\infty\}$

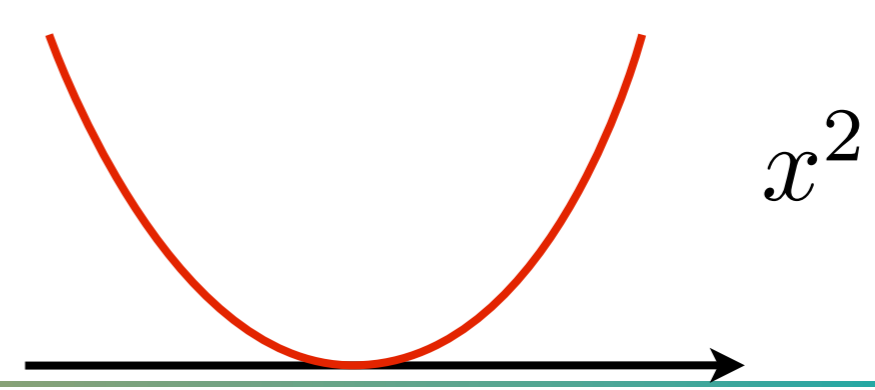
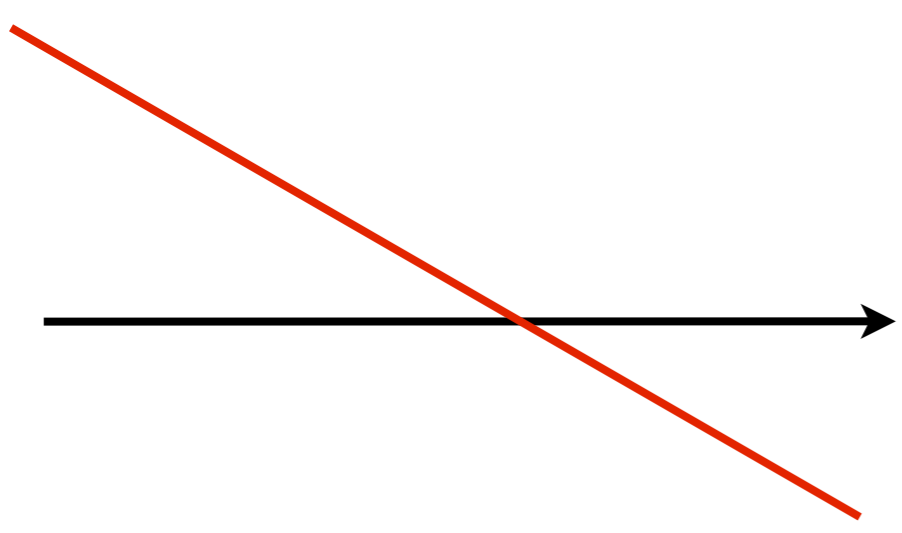
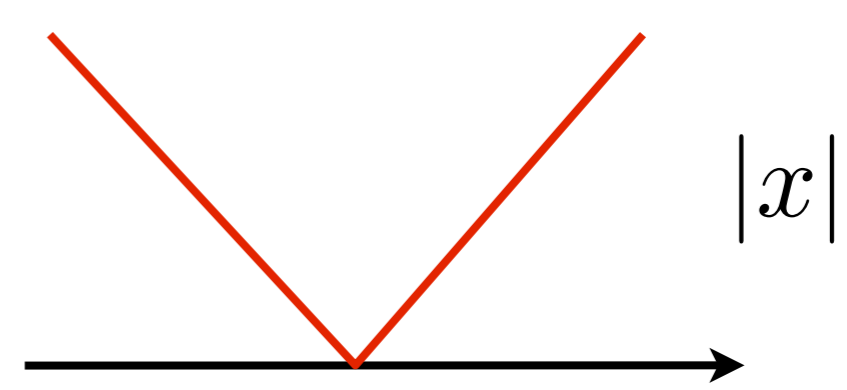
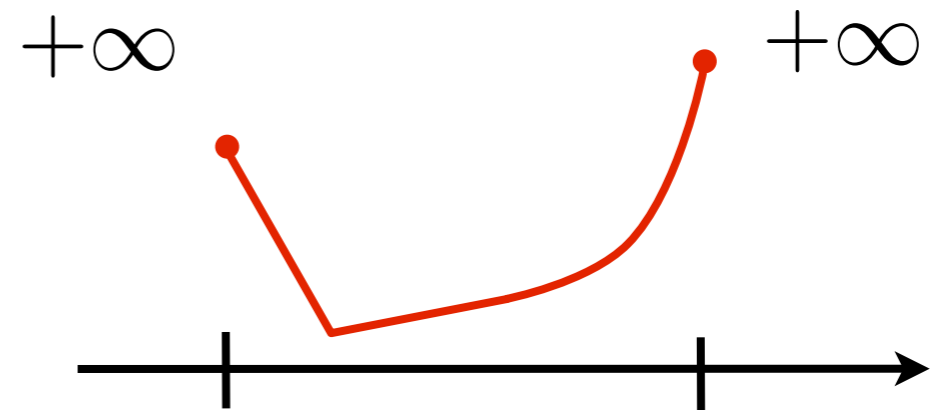
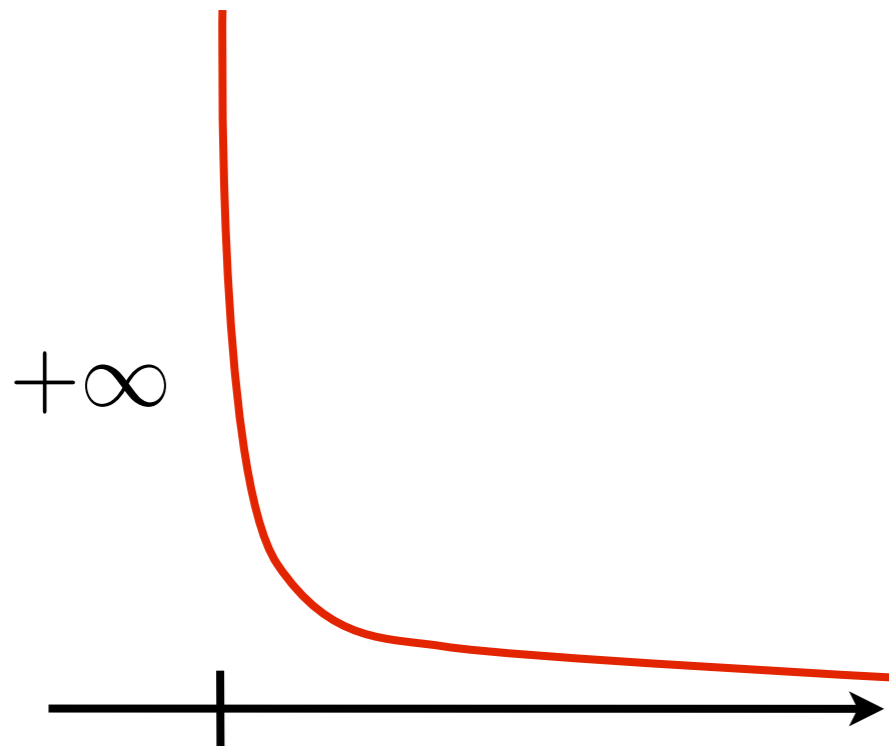
Convex functions

f is **convex** if $\forall x, y \in \mathcal{X}$ and $a \in [0, 1]$,
 $f(ax + (1 - a)y) \leq af(x) + (1 - a)f(y)$



Convex functions

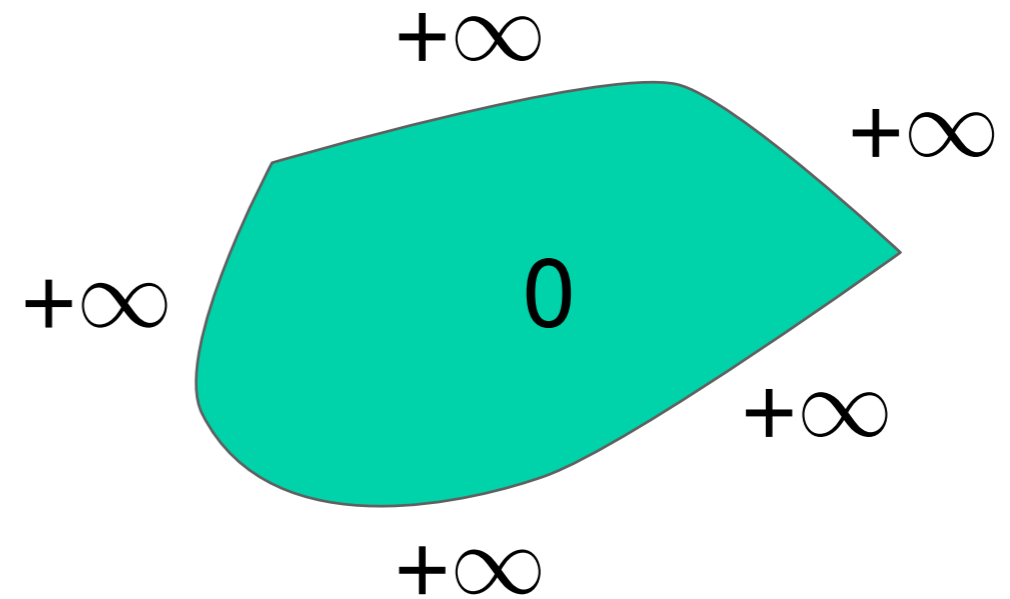
Some convex functions: $\mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$



Indicator functions

For a closed convex set $\Omega \subset \mathcal{X}$, its **indicator function** is

$$I_{\Omega}(x) = \begin{cases} 0 & \text{if } x \in \Omega, \\ +\infty & \text{else.} \end{cases}$$

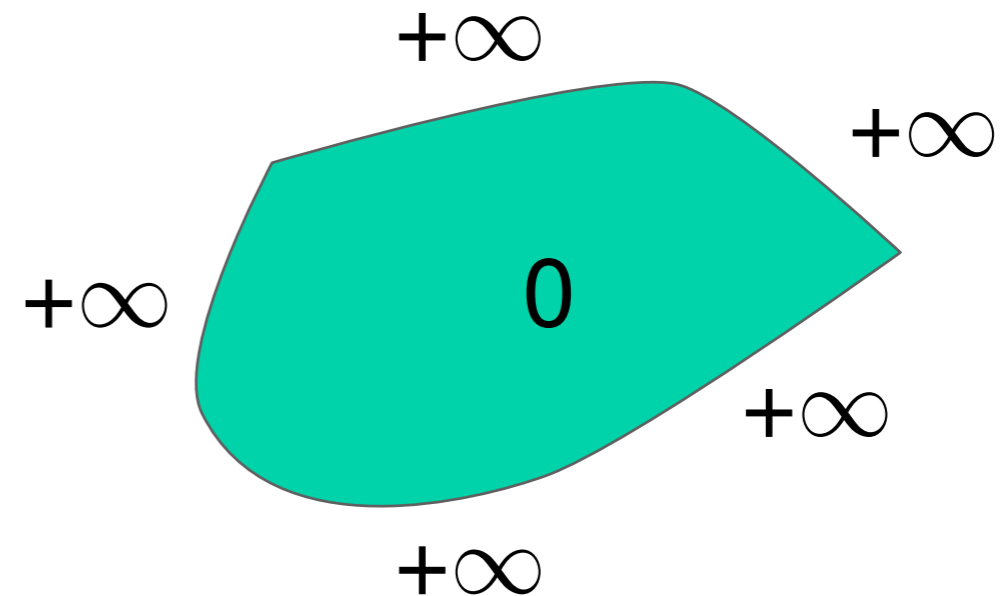


I_{Ω} is convex.

Indicator functions

For a closed convex set $\Omega \subset \mathcal{X}$, its **indicator function** is

$$I_{\Omega}(x) = \begin{cases} 0 & \text{if } x \in \Omega, \\ +\infty & \text{else.} \end{cases}$$



I_{Ω} is convex.

Note: $I_{\Omega_1} + I_{\Omega_2} = I_{\Omega_1 \cap \Omega_2}$

Constrained optimization

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \Psi(x) = \sum_{m=1}^M g_m(L_m x)$$

encompasses the presence of constraints:

$$\begin{aligned} \underset{x \in \Omega}{\text{minimize}} \ f(x) &\equiv \underset{x \in \mathcal{X}}{\text{minimize}} \ f(x) \ \text{s.t. } x \in \Omega \\ &\equiv \underset{x \in \mathcal{X}}{\text{minimize}} \ f(x) + I_{\Omega}(x) \end{aligned}$$

Optimization algorithms

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \Psi(x) = \sum_{m=1}^M g_m(L_m x)$$



iterative algorithm computing

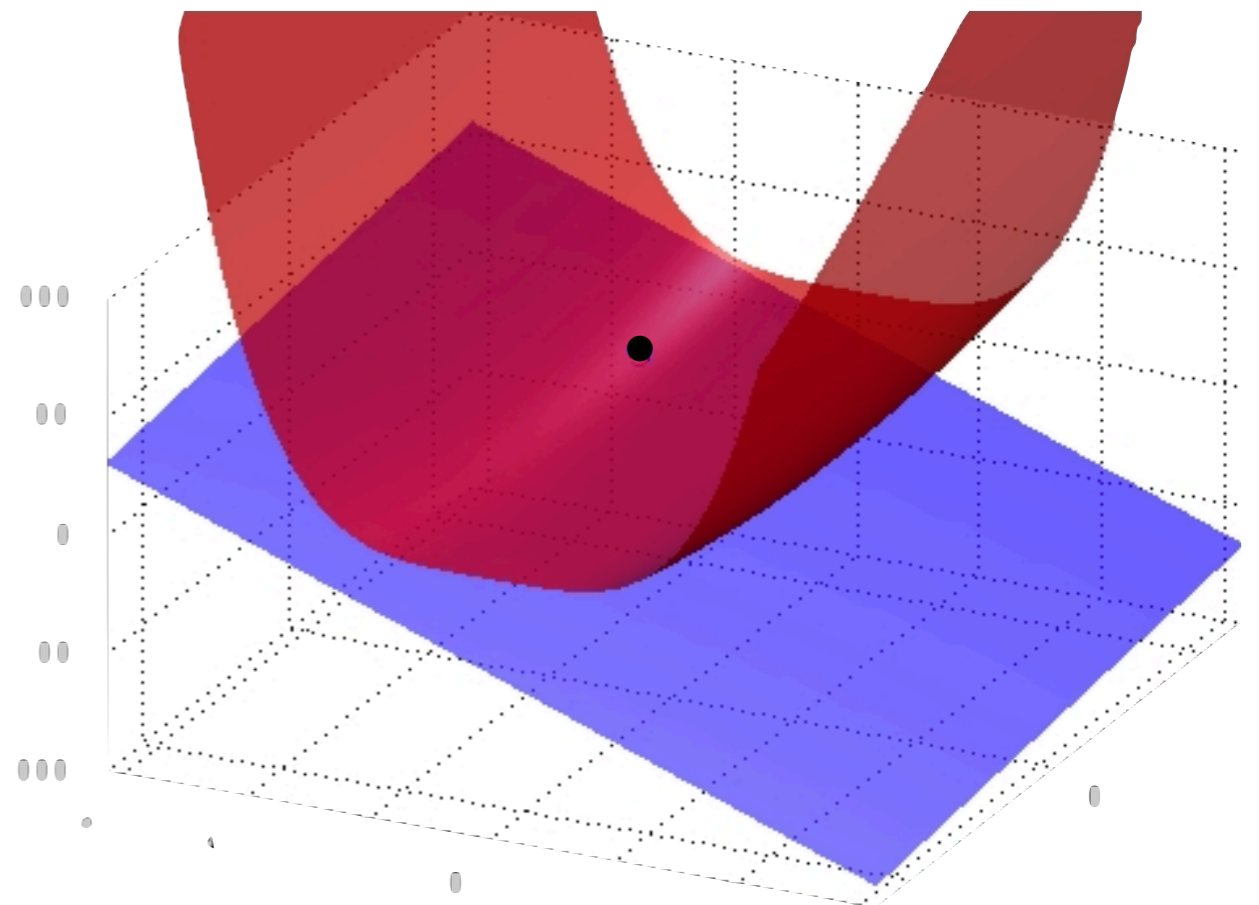
$$x^{k+1} = T(x^k)$$

with x^k converging to $x^* = T(x^*)$

The gradient

$f : \mathcal{X} \rightarrow \mathbb{R}$ is differentiable (=smooth) at x
if there exists a unique element $\nabla f(x) \in \mathcal{X}$ such that

$$\forall \mathbf{e} \in \mathcal{X}, f(x + \mathbf{e}) = f(x) + \langle \mathbf{e}, \nabla f(x) \rangle + o(\|\mathbf{e}\|)$$




The gradient

$f : \mathcal{X} \rightarrow \mathbb{R}$ is differentiable (=smooth) at x
if there exists a unique element $\nabla f(x) \in \mathcal{X}$ such that

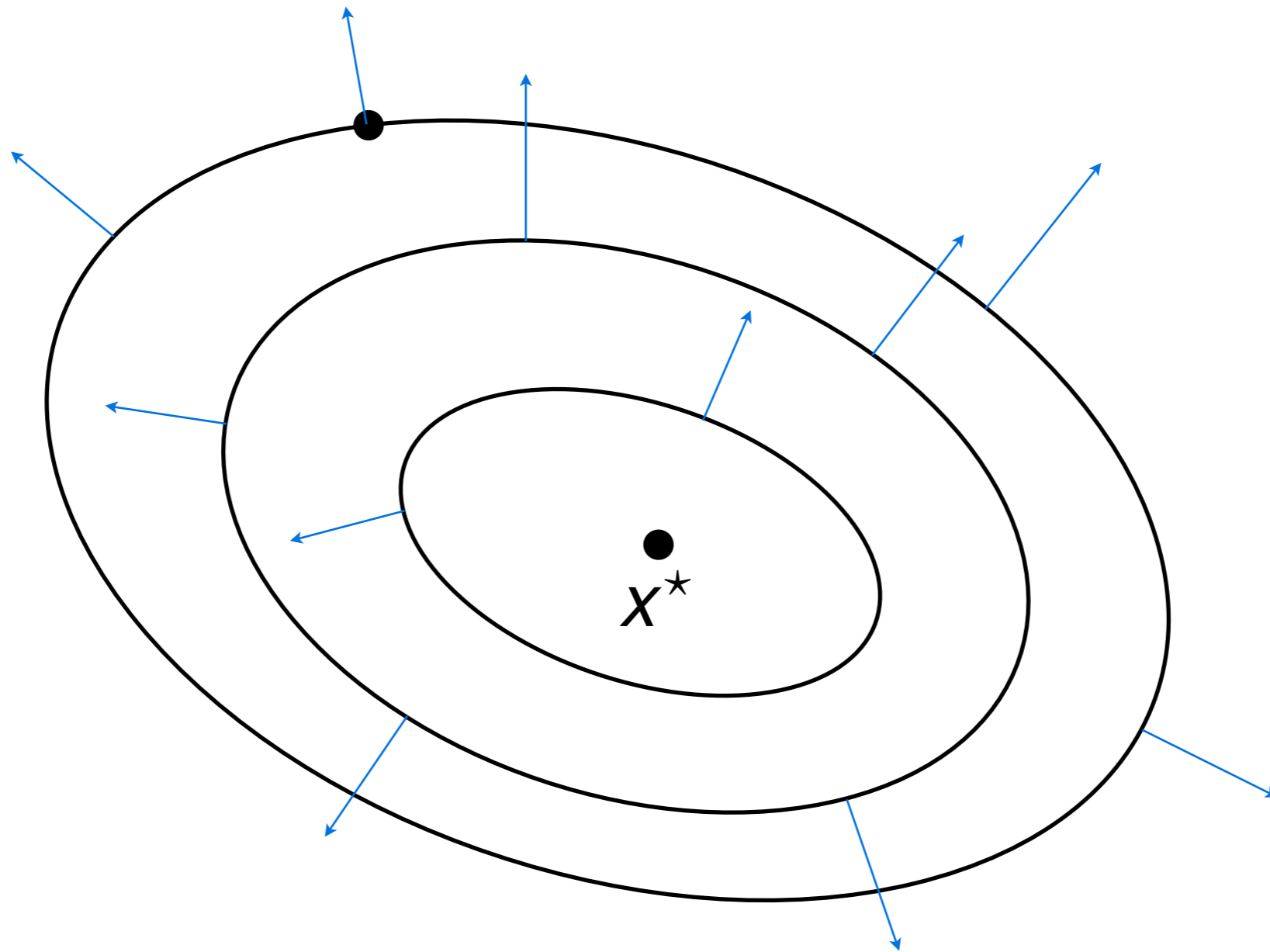
$$\forall \mathbf{e} \in \mathcal{X}, f(x + \mathbf{e}) = f(x) + \langle \mathbf{e}, \nabla f(x) \rangle + o(\|\mathbf{e}\|)$$

$$\psi(x) = \sum_{m=1}^M g_m(L_m x)$$

 $\nabla \psi(x) = \sum_{m=1}^M L_m^* \nabla g_m(L_m x)$

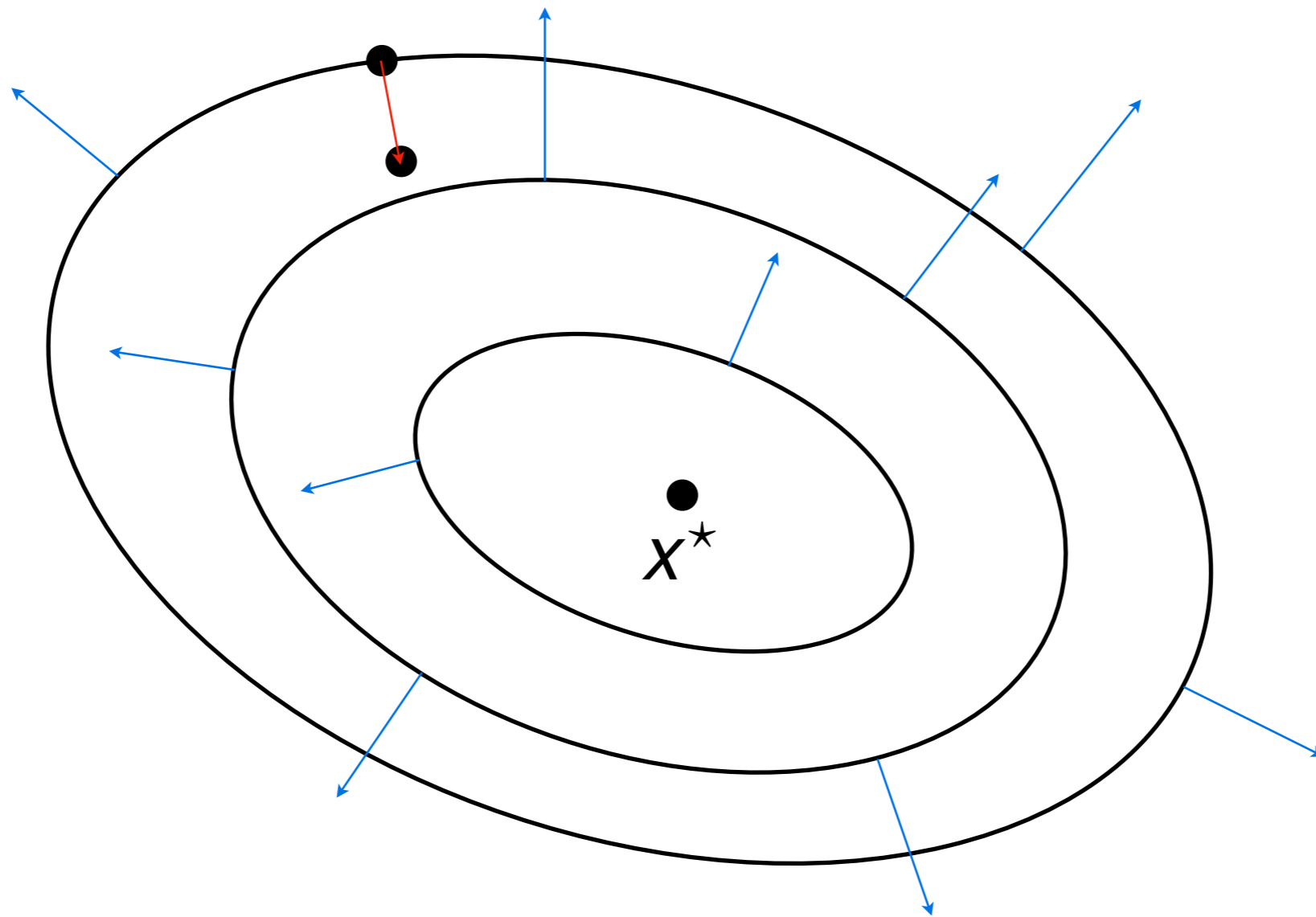
Gradient descent

$$x^{k+1} = x^k - \gamma \nabla \Psi(x^k)$$



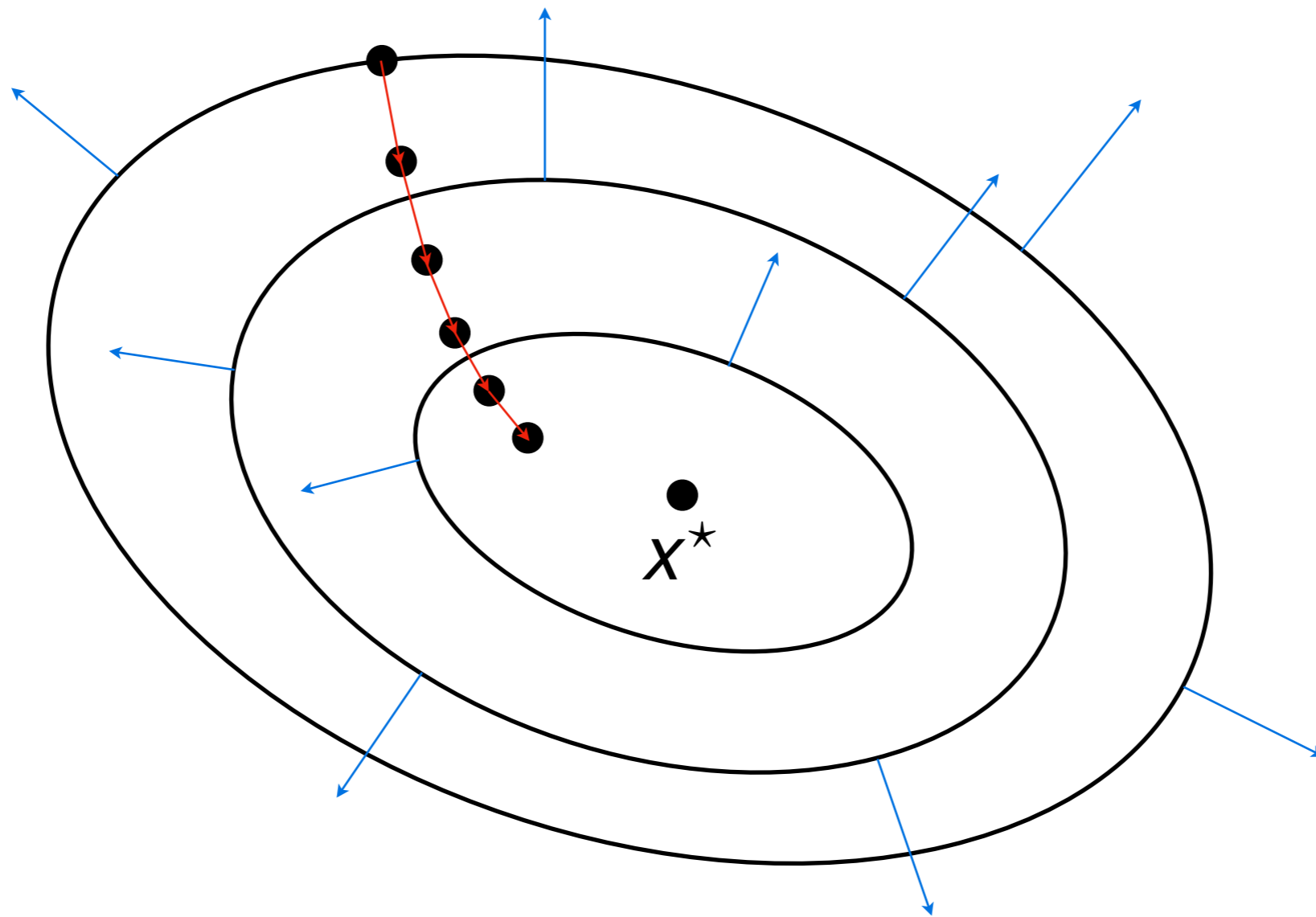
Gradient descent

$$x^{k+1} = x^k - \gamma \nabla \Psi(x^k)$$

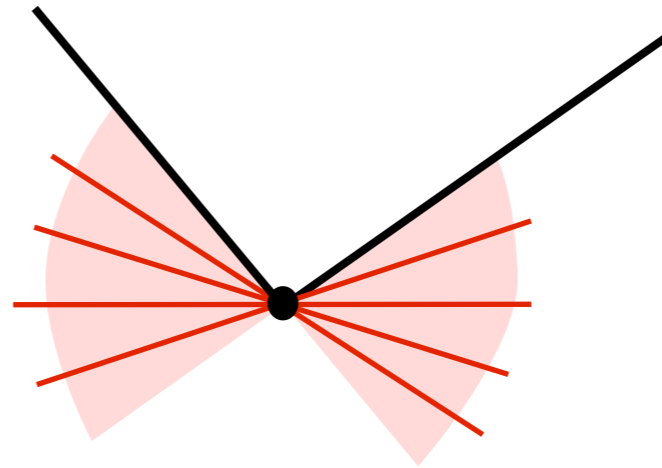


Gradient descent

$$x^{k+1} = x^k - \gamma \nabla \Psi(x^k)$$



The subdifferential



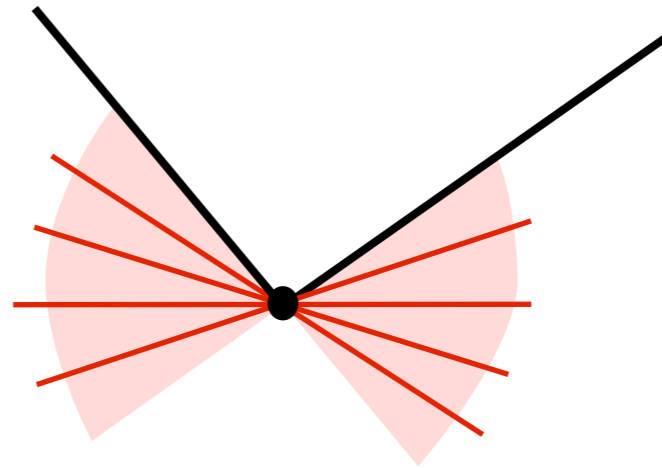
$$\partial f: \mathcal{X} \rightarrow 2^{\mathcal{X}}$$

$$x \mapsto \{u \in \mathcal{X} : \forall y \in \mathcal{X}, f(x) + \langle y - x, u \rangle \leq f(y)\}$$



$\partial f(x)$ is the set of gradients of the affine minorants of f at x

The subdifferential

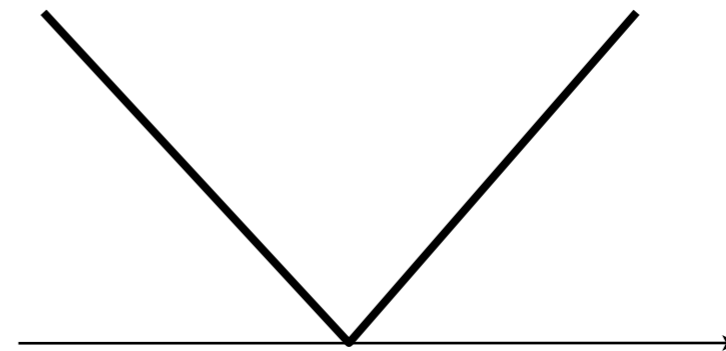


$$\partial f: \mathcal{X} \rightarrow 2^{\mathcal{X}}$$

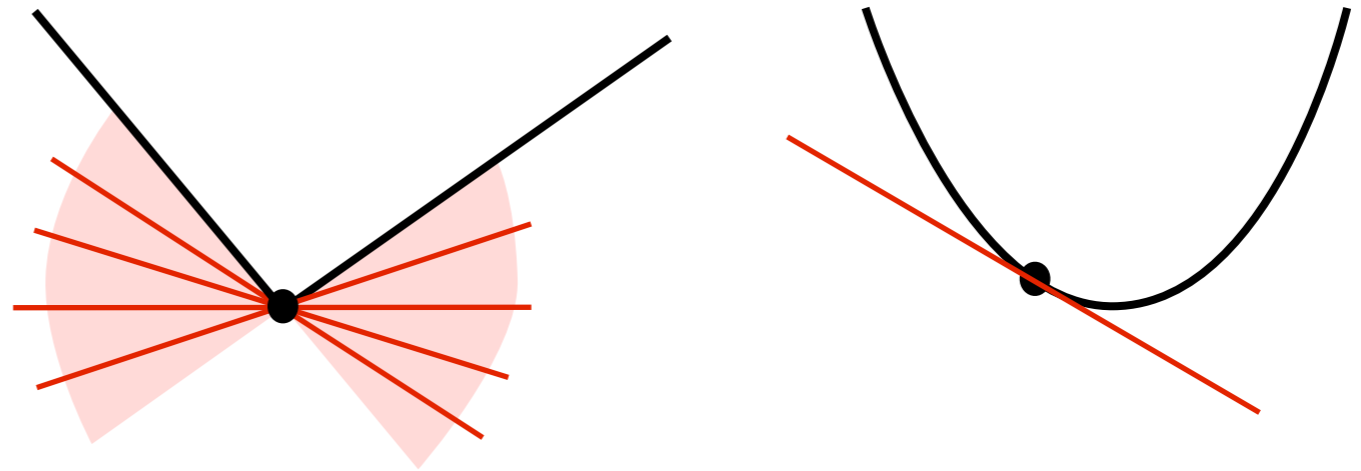
$$x \mapsto \{u \in \mathcal{X} : \forall y \in \mathcal{X}, f(x) + \langle y - x, u \rangle \leq f(y)\}$$

Example: $f = |\cdot|$

👉 $\partial f(0) = [-1, 1]$



The subdifferential



$$\partial f: \mathcal{X} \rightarrow 2^{\mathcal{X}}$$

$$x \mapsto \{u \in \mathcal{X} : \forall y \in \mathcal{X}, f(x) + \langle y - x, u \rangle \leq f(y)\}$$

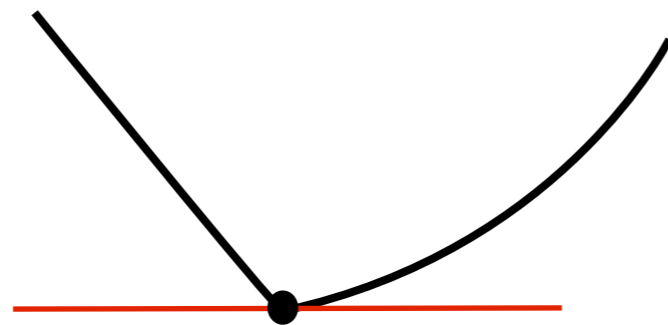
f is convex and smooth at $x \rightarrow \partial f(x) = \{\nabla f(x)\}$.

Fermat's rule

$$x^* \in \arg \min f$$

\Leftrightarrow

$$0 \in \partial f(x^*)$$



Pierre de Fermat,
1601-1665

First-order conditions

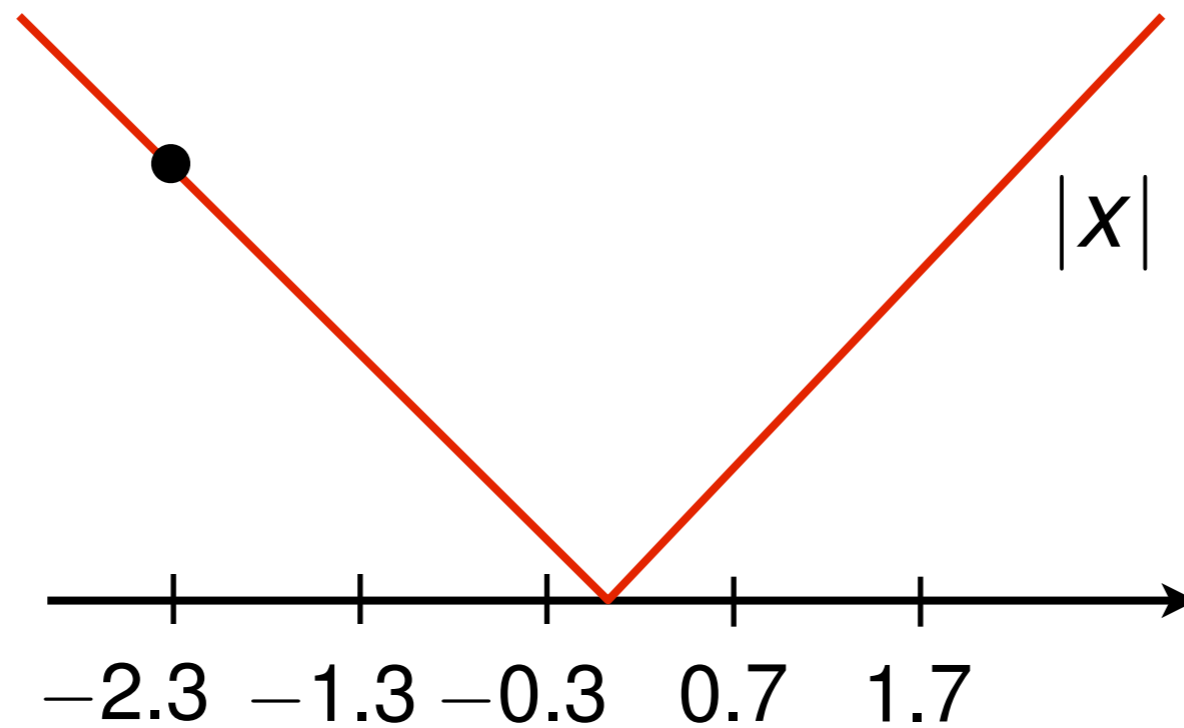
$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \psi(x) = \sum_{m=1}^M g_m(L_m x)$$

↑

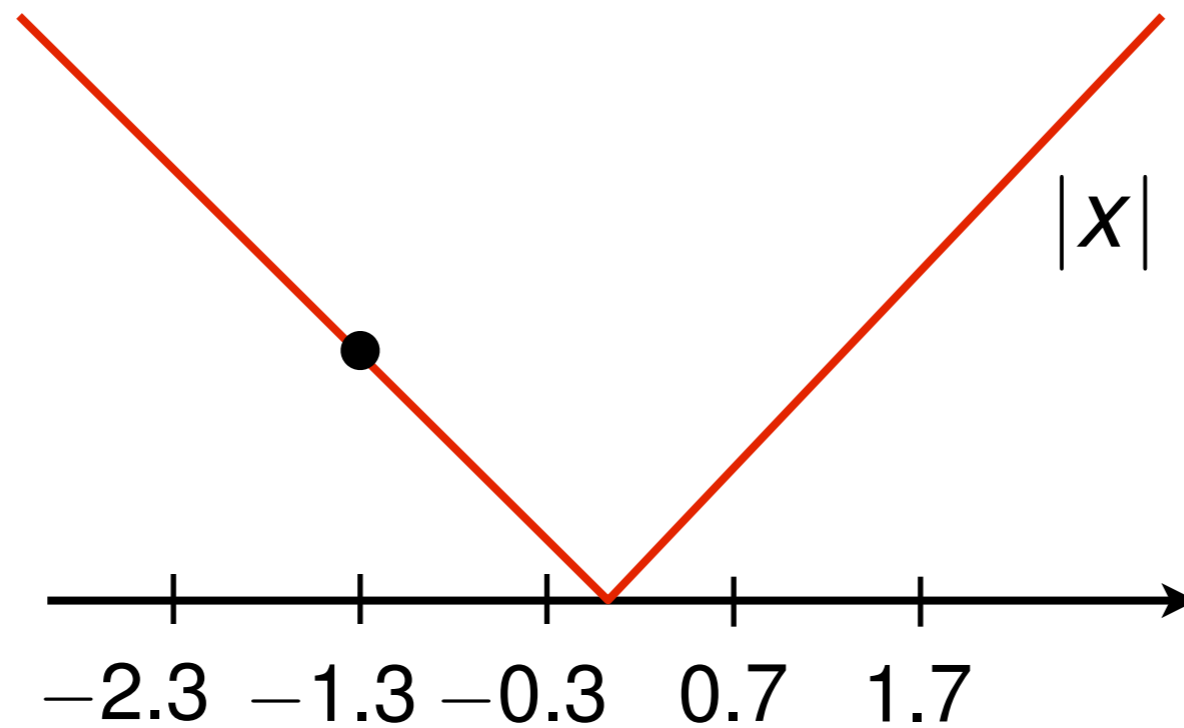
find $x^* \in \mathcal{X}$ such that

$$0 \in \sum_{m=1}^M L_m^* \partial g_m(L_m x^*)$$

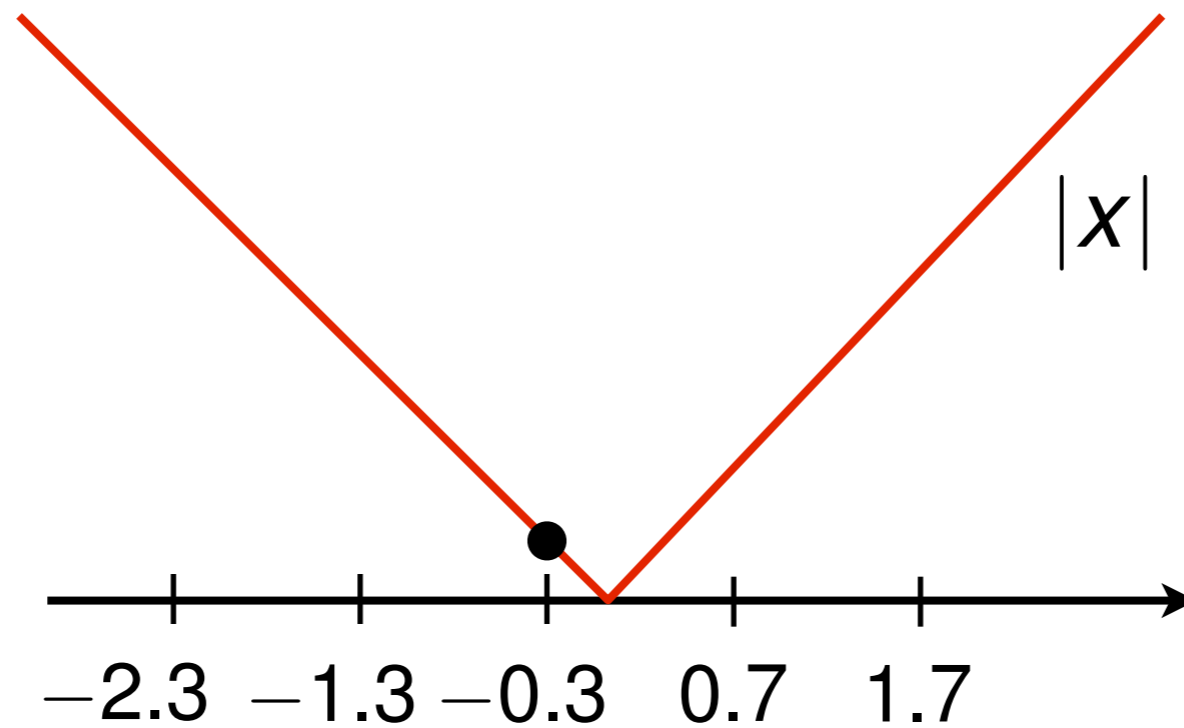
Nonsmooth minimization?



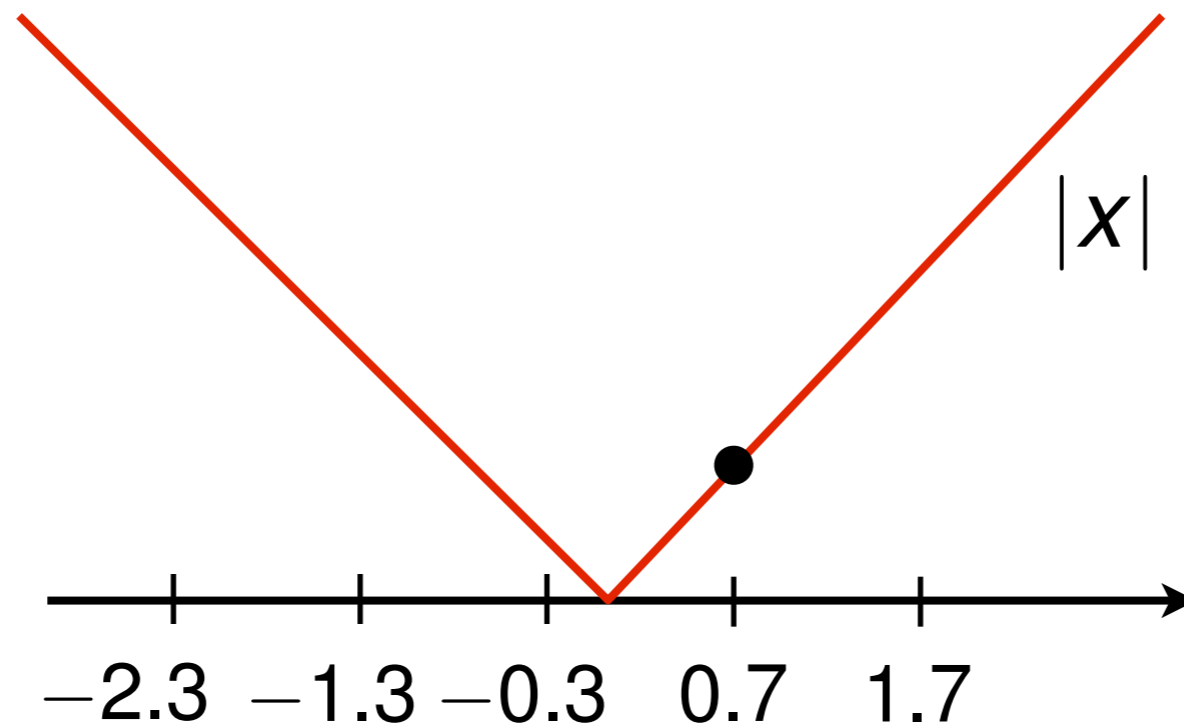
Nonsmooth minimization?



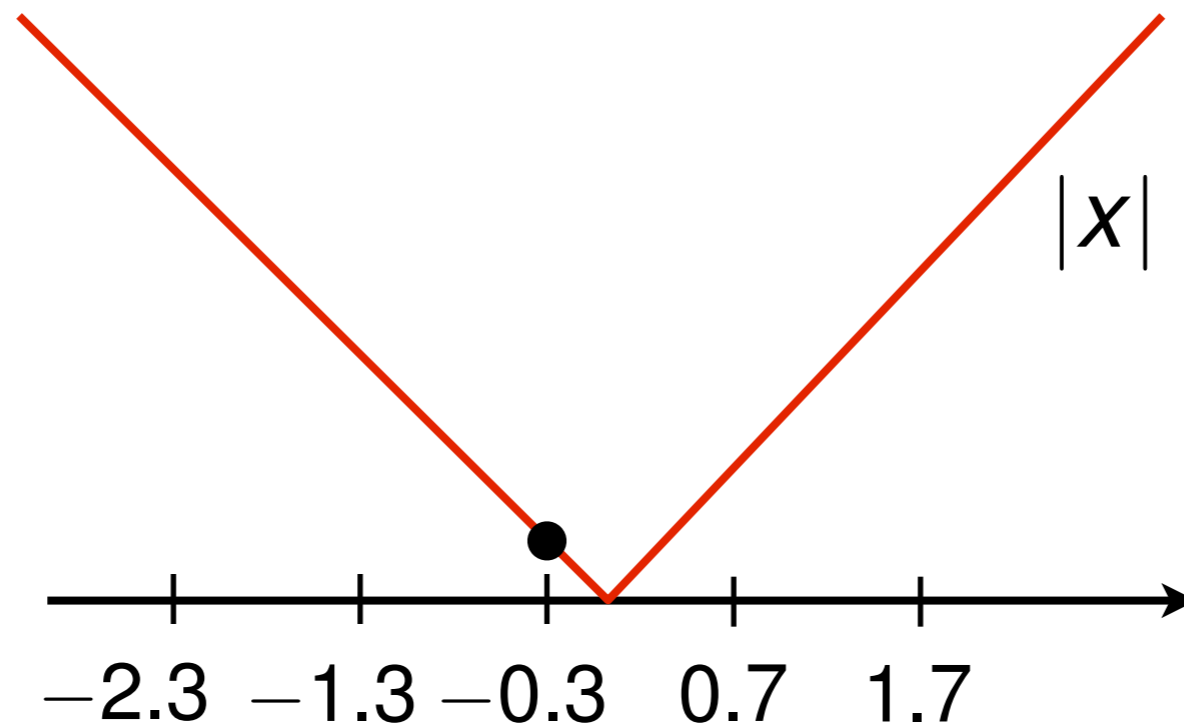
Nonsmooth minimization?



Nonsmooth minimization?

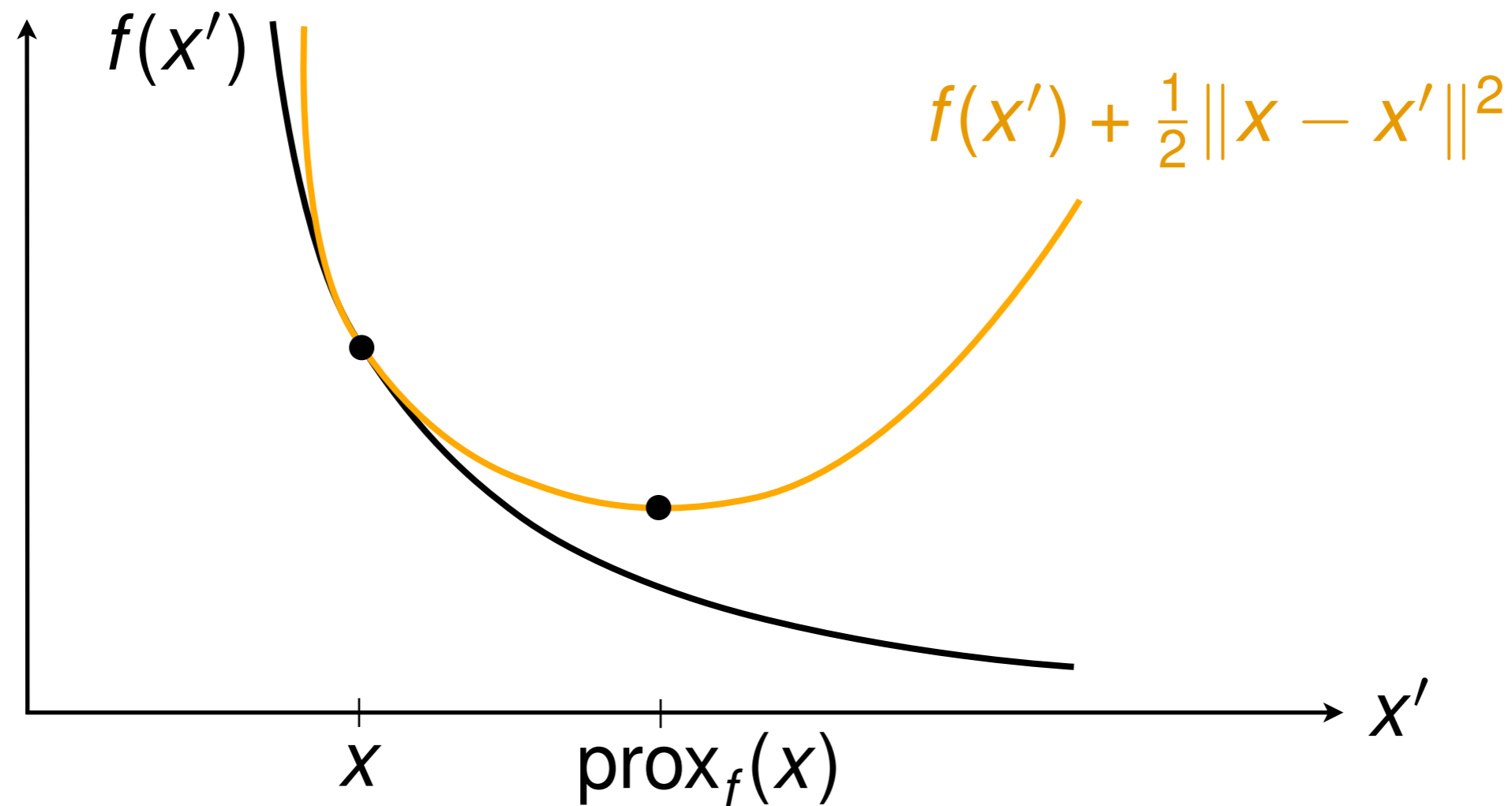


Nonsmooth minimization?



The proximity operator

$$\text{prox}_f : \mathcal{X} \rightarrow \mathcal{X} : x \mapsto \arg \min_{x' \in \mathcal{X}} \left(f(x') + \frac{1}{2} \|x - x'\|^2 \right)$$



The proximity operator

$$\text{prox}_f : \mathcal{X} \rightarrow \mathcal{X} : x \mapsto \arg \min_{x' \in \mathcal{X}} \left(f(x') + \frac{1}{2} \|x - x'\|^2 \right)$$

$$\text{prox}_f = (\partial f + \text{Id})^{-1}$$

The proximity operator

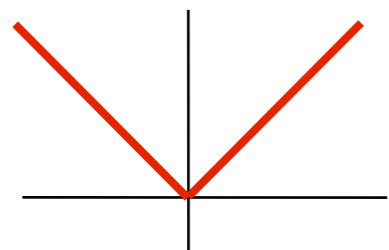
$$\text{prox}_f : \mathcal{X} \rightarrow \mathcal{X} : x \mapsto \arg \min_{x' \in \mathcal{X}} \left(f(x') + \frac{1}{2} \|x - x'\|^2 \right)$$

$$\text{prox}_f = (\partial f + \text{Id})^{-1}$$

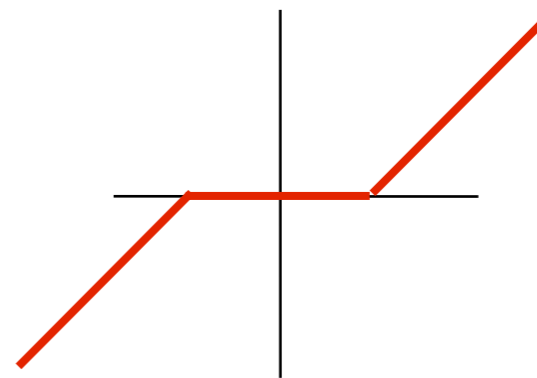
$$y = \text{prox}_f(x) \quad \equiv \quad y \in x - \partial f(y)$$

The proximity operator

$$\text{prox}_f : \mathcal{X} \rightarrow \mathcal{X} : x \mapsto \arg \min_{x' \in \mathcal{X}} \left(f(x') + \frac{1}{2} \|x - x'\|^2 \right)$$



$$f(x) = |x|$$

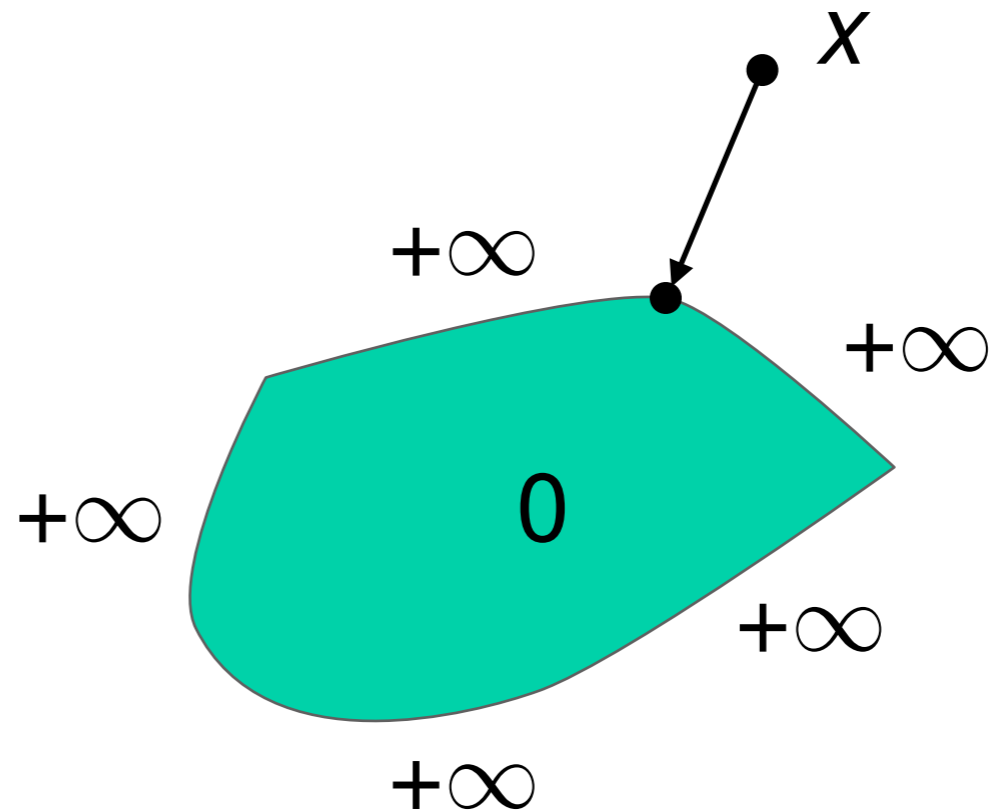


$$\text{prox}_f(x) = \text{sgn}(x) \max(|x| - 1, 0)$$

The proximity operator

$$\text{prox}_f : \mathcal{X} \rightarrow \mathcal{X} : x \mapsto \arg \min_{x' \in \mathcal{X}} \left(f(x') + \frac{1}{2} \|x - x'\|^2 \right)$$

$$\text{prox}_{I_\Omega} = \text{proj}_\Omega$$



The proximity operator

Exact, finite time, algorithms are available to compute the proximity operator of:

- $\|X\|_* \rightarrow \text{SVD}, O(N^3)$
- 1-D TV \rightarrow taut-string alg., $O(N)$
- 2-D anisotropic TV \rightarrow graph cuts
- proj. onto the simplex $\rightarrow O(N)$

...

The proximity operator

Exact, finite time, algorithms are available to compute the proximity operator of:

- $\|X\|_* \rightarrow \text{SVD}, O(N^3)$
- 1-D TV \rightarrow taut-string alg., $O(N)$
- 2-D anisotropic TV \rightarrow graph cuts
- proj. onto the simplex $\rightarrow O(N)$

...

LC, "A direct algorithm for 1-D total variation...", 2013

LC, "Fast projection onto the simplex...", 2016

Pustelnik, LC, "Proximity operator of a sum...", 2017

Proximal splitting algorithms

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \sum_{m=1}^M g_m(L_m x)$$



No easy form of $\text{prox}_{g_1+g_2}$ or $\text{prox}_{g \circ L}$

Proximal splitting algorithms

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \sum_{m=1}^M g_m(L_m x)$$



We want **full splitting**, with individual activation of L_m , L_m^* , the gradient or proximity operator of g_m .

Proximal splitting algorithms

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \left(f(x) + \sum_{m=1}^M g_m(L_m x) + h(x) \right)$$

with:

- h smooth with β -Lipschitz continuous gradient
→ calls to ∇h
- simple functions f and g_m → calls to $\text{prox}_{\gamma f}$ and $\text{prox}_{\tau_m g_m}$
- calls to L_m, L_m^*

Product space trick

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \left(f(x) + g(Lx) + h(x) \right)$$

$$g(u) = \sum_{m=1}^M g_m(u_m)$$



$$g(Lx) = \sum_{m=1}^M g_m(L_m x)$$

$$Lx = (L_1 x, \dots, L_m x)$$

Minimization of 3 functions

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \left(\underset{\downarrow}{f(x)} + \underset{\downarrow}{g(Lx)} + \underset{\downarrow}{h(x)} \right)$$

$\text{prox}_{\gamma f}, \quad \text{prox}_{\tau g}, \quad \nabla h, \quad L, \quad L^*$

Proximal splitting algorithms

minimize $f + g \circ L + h$

$f + h$



forward-backward

$f + g$






Douglas-Rachford / ADMM

1979



Proximal splitting algorithms

minimize $f + g \circ L + h$

	$f + h$		forward-backward
1979	$f + g$		Douglas-Rachford / ADMM
2011	$f + g \circ L$		Chambolle-Pock






Proximal splitting algorithms

minimize $f + g \circ L + h$

	$f + h$		forward-backward
1979	$f + g$		Douglas-Rachford / ADMM
2011	$f + g \circ L$		Chambolle-Pock
2011	$g \circ L + h$		PAPC

Proximal splitting algorithms

minimize $f + g \circ L + h$







	$f + h$		forward-backward
1979	$f + g$		Douglas-Rachford / ADMM
2011	$f + g \circ L$		Chambolle-Pock
2011	$g \circ L + h$		PAPC
2013	$f + g \circ L + h$		Condat, Vu

LC, "A primal-dual splitting method for convex optimization...", 2013

Vu, "A splitting algorithm for dual monotone inclusions...", 2013








Proximal splitting algorithms

minimize $f + g \circ L + h$

	$f + h$		forward-backward
1979	$f + g$		Douglas-Rachford / ADMM
2011	$f + g \circ L$		Chambolle-Pock
2011	$g \circ L + h$		PAPC
2013	$f + g \circ L + h$		Condat, Vu
2017	$f + g + h$		Davis-Yin









Proximal splitting algorithms

minimize $f + g \circ L + h$

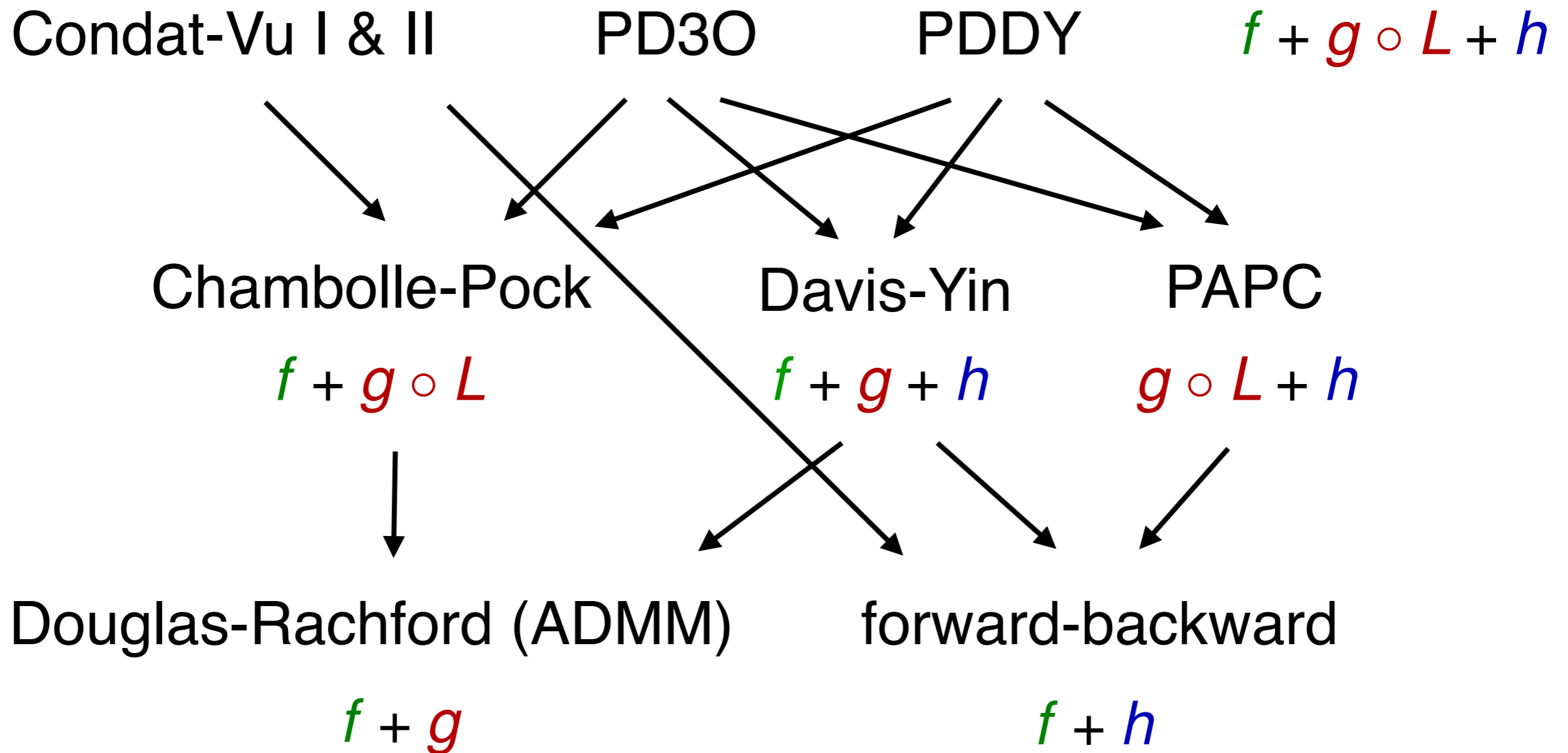
	$f + h$		forward-backward
1979	$f + g$		Douglas-Rachford / ADMM
2011	$f + g \circ L$		Chambolle-Pock
2011	$g \circ L + h$		PAPC
2013	$f + g \circ L + h$		Condat, Vu
2017	$f + g + h$		Davis-Yin
2018	$f + g \circ L + h$		PD3O

Proximal splitting algorithms

minimize $f + g \circ L + h$

	$f + h$		forward-backward	
1979	$f + g$		Douglas-Rachford / ADMM	
2011	$f + g \circ L$		Chambolle-Pock	
2011	$g \circ L + h$		PAPC	
2013	$f + g \circ L + h$		Condat, Vu	
2017	$f + g + h$		Davis-Yin	
2018	$f + g \circ L + h$		PD3O	
2022	$f + g \circ L + h$		PDDY	Salim, LC et al., "Dualize, split, randomize...", 2022

Proximal splitting algorithms



4 primal-dual algorithms

Condat–Vu algorithm form I

$$\begin{cases} x^{k+1} = \text{prox}_{\gamma f} (x^k - \gamma \nabla h(x^k) - \gamma L^* u^k) \\ u^{k+1} = \text{prox}_{\tau g^*} (u^k + \tau L(2x^{k+1} - x^k)) \end{cases}$$

minimize
 $f + g \circ L + h$

Condat–Vu algorithm form II

$$\begin{cases} u^{k+1} = \text{prox}_{\tau g^*} (u^k + \tau Lx^k) \\ x^{k+1} = \text{prox}_{\gamma f} (x^k - \gamma \nabla h(x^k) - \gamma L^*(2u^{k+1} - u^k)) \end{cases}$$

PD3O algorithm

$$\begin{cases} x^{k+1} = \text{prox}_{\gamma f} (x^k - \gamma \nabla h(x^k) - \gamma L^* u^k) \\ u^{k+1} = \text{prox}_{\tau g^*} (u^k + \tau L(2x^{k+1} - x^k - \gamma \nabla h(x^{k+1}) + \gamma \nabla h(x^k))) \end{cases}$$

PDDY algorithm

$$\begin{cases} u^{k+1} = \text{prox}_{\tau g^*} (u^k + \tau Lx^k) \\ x^{k+1} = \text{prox}_{\gamma f} (x^k - \gamma \nabla h(x^k - \gamma L^*(u^{k+1} - u^k)) - \gamma L^*(2u^{k+1} - u^k)) \end{cases}$$

First-order conditions

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \left(f(x) + g(Lx) + h(x) \right)$$

$\sim \equiv$

$$0 \in \nabla h(x^*) + \partial f(x^*) + L^* \partial g(Lx^*)$$

First-order conditions

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \left(f(x) + g(Lx) + h(x) \right)$$

$\sim \equiv$

$$0 \in \nabla h(x^*) + \partial f(x^*) + L^* \partial g(Lx^*)$$

\equiv

$$\begin{cases} 0 \in \nabla h(x^*) + \partial f(x^*) + L^* u^* \\ u^* \in \partial g(Lx^*) \end{cases}$$

First-order conditions

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \left(f(x) + g(Lx) + h(x) \right)$$

$\sim \equiv$

$$0 \in \nabla h(x^*) + \partial f(x^*) + L^* \partial g(Lx^*)$$

\equiv

$$\begin{cases} 0 \in \nabla h(x^*) + \partial f(x^*) + L^* u^* \\ 0 \in -Lx^* + \partial g^*(u^*) \end{cases}$$

$$(\partial g^* = (\partial g)^{-1})$$

First-order conditions

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \left(f(x) + g(Lx) + h(x) \right)$$

$\sim \equiv$

Find (x^*, u^*) solution to

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} \nabla h(x) + \partial f(x) + L^* u \\ -Lx + \partial g^*(u) \end{pmatrix}$$

First-order conditions

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \left(f(x) + g(Lx) + h(x) \right)$$

$\sim \equiv$

Find (x^*, u^*) solution to

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} \nabla h(x) + \partial f(x) + L^* u \\ -Lx + \partial g^*(u) \end{pmatrix}$$

Dual problem:

$$\text{Find } u^* \in \arg \min_{u \in \mathcal{U}} \left((f + h)^*(-L^* u) + g^*(u) \right)$$

First-order conditions

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \left(f(x) + g(Lx) + h(x) \right)$$

$\sim \equiv$

Find (x^*, u^*) solution to

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} \nabla h(x) + \partial f(x) + L^* u \\ -Lx + \partial g^*(u) \end{pmatrix}$$



monotone

M monotone: $\forall (x, x') \in \mathcal{X}^2, v \in Mx, v' \in Mx', \langle x - x', v - v' \rangle \geq 0$

Primal-dual 3 operator splitting

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} L^* u \\ -Lx + \partial g^*(u) \end{pmatrix} + \begin{pmatrix} \partial f(x) \\ 0 \end{pmatrix} + \begin{pmatrix} \nabla h(x) \\ 0 \end{pmatrix}$$

Primal-dual 3 operator splitting

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} L^* u \\ -Lx + \partial g^*(u) \end{pmatrix} + \begin{pmatrix} \partial f(x) \\ 0 \end{pmatrix} + \begin{pmatrix} \nabla h(x) \\ 0 \end{pmatrix}$$

or


$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} \partial f(x) + L^* u \\ -Lx \end{pmatrix} + \begin{pmatrix} 0 \\ \partial g^*(u) \end{pmatrix} + \begin{pmatrix} \nabla h(x) \\ 0 \end{pmatrix}$$


Primal-dual 3 operator splitting

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} L^*u \\ -Lx + \partial g^*(u) \end{pmatrix} + \begin{pmatrix} \partial f(x) \\ 0 \end{pmatrix} + \begin{pmatrix} \nabla h(x) \\ 0 \end{pmatrix}$$

or

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} \partial f(x) + L^*u \\ -Lx \end{pmatrix} + \begin{pmatrix} 0 \\ \partial g^*(u) \end{pmatrix} + \begin{pmatrix} \nabla h(x) \\ 0 \end{pmatrix}$$


 monotone


 cocoercive

C is ξ -cocoercive: $\forall (x, x') \in \mathcal{X}^2, \xi \|Cx - Cx'\|^2 \leq \langle x - x', Cx - Cx' \rangle$

Primal-dual 3 operator splitting

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} L^* u \\ -Lx + \partial g^*(u) \end{pmatrix} + \begin{pmatrix} \partial f(x) \\ 0 \end{pmatrix} + \begin{pmatrix} \nabla h(x) \\ 0 \end{pmatrix}$$

or

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} \partial f(x) + L^* u \\ -Lx \end{pmatrix} + \begin{pmatrix} 0 \\ \partial g^*(u) \end{pmatrix} + \begin{pmatrix} \nabla h(x) \\ 0 \end{pmatrix}$$



$$0 \in A(z^*) + B(z^*) + C(z^*)$$

Primal-dual 3 operator splitting

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} L^* u \\ -Lx + \partial g^*(u) \end{pmatrix} + \begin{pmatrix} \partial f(x) \\ 0 \end{pmatrix} + \begin{pmatrix} \nabla h(x) \\ 0 \end{pmatrix}$$

or

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} \partial f(x) + L^* u \\ -Lx \end{pmatrix} + \begin{pmatrix} 0 \\ \partial g^*(u) \end{pmatrix} + \begin{pmatrix} \nabla h(x) \\ 0 \end{pmatrix}$$



$$0 \in P^{-1} A(z^*) + P^{-1} B(z^*) + P^{-1} C(z^*)$$

for some $P \succ 0$

(M monotone $\Rightarrow P^{-1} M$ monotone w.r.t. $\langle \cdot, P \cdot \rangle$)

Primal-dual 3 operator splitting

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} L^* u \\ -Lx + \partial g^*(u) \end{pmatrix} + \begin{pmatrix} \partial f(x) \\ 0 \end{pmatrix} + \begin{pmatrix} \nabla h(x) \\ 0 \end{pmatrix}$$

or

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} \partial f(x) + L^* u \\ -Lx \end{pmatrix} + \begin{pmatrix} 0 \\ \partial g^*(u) \end{pmatrix} + \begin{pmatrix} \nabla h(x) \\ 0 \end{pmatrix}$$

$$0 \in P^{-1} A(z^*) + P^{-1} B(z^*) + P^{-1} C(z^*)$$

Davis-Yin splitting algorithm

$$\begin{cases} z^k = J_{P^{-1}B}(v^k) \\ w^{k+1} = J_{P^{-1}A}(2z^k - v^k - P^{-1}C(z^k)) \\ v^{k+1} = v^k + w^{k+1} - z^k \end{cases} \quad (J_M = (M + \text{Id})^{-1})$$

PD3O algorithm

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} L^* u \\ -Lx + \partial g^*(u) \end{pmatrix} + \begin{pmatrix} \partial f(x) \\ 0 \end{pmatrix} + \begin{pmatrix} \nabla h(x) \\ 0 \end{pmatrix}$$

👉 Davis-Yin splitting with $P = \begin{pmatrix} \frac{1}{\gamma} \text{Id} & 0 \\ 0 & \frac{1}{\tau} \text{Id} - \gamma LL^* \end{pmatrix}$

$$\equiv$$

PD3O algorithm

$$\begin{cases} x^{k+1} = \text{prox}_{\gamma f}(s^k) \\ u^{k+1} = \text{prox}_{\tau g^*} \left(u^k + \tau L(2x^{k+1} - s^k - \gamma \nabla h(x^{k+1}) - \gamma L^* u^k) \right) \\ s^{k+1} = x^{k+1} - \gamma \nabla h(x^{k+1}) - \gamma L^* u^{k+1} \end{cases}$$

PDDY algorithm

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} L^* u \\ -Lx + \partial g^*(u) \end{pmatrix} + \begin{pmatrix} \partial f(x) \\ 0 \end{pmatrix} + \begin{pmatrix} \nabla h(x) \\ 0 \end{pmatrix}$$



Davis-Yin splitting with
and swapping A and B

$$P = \begin{pmatrix} \frac{1}{\gamma} \text{Id} & 0 \\ 0 & \frac{1}{\tau} \text{Id} - \gamma LL^* \end{pmatrix}$$


$$\equiv$$

PDDY algorithm

$$\begin{cases} \hat{x}^k = \text{prox}_{\gamma f}(x^k - \gamma \nabla h(x^k) + \gamma h^k) \\ u^{k+1} = \text{prox}_{\tau g^*}(u^k + \tau L \hat{x}^k) \\ h^{k+1} = -L^* u^{k+1} \\ x^{k+1} = \hat{x}^k + \gamma(h^{k+1} - h^k) \end{cases}$$

Condat-Vu algorithm

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} \partial f(x) + L^* u \\ -Lx \end{pmatrix} + \begin{pmatrix} 0 \\ \partial g^*(u) \end{pmatrix} + \begin{pmatrix} \nabla h(x) \\ 0 \end{pmatrix}$$

 Davis-Yin splitting with $P = \begin{pmatrix} \frac{1}{\gamma} \text{Id} & -\tau L^* L & 0 \\ 0 & 0 & \frac{1}{\tau} \text{Id} \end{pmatrix}$

\equiv

Condat-Vu algorithm

$$\begin{cases} x^{k+1} = \text{prox}_{\gamma f} (x^k - \gamma \nabla h(x^k) - \gamma L^* u^k) \\ u^{k+1} = \text{prox}_{\tau g^*} (u^k + \tau Lx(2x^{k+1} - x^k)) \end{cases}$$

Condat-Vu algorithm

$$\begin{pmatrix} 0 \\ 0 \end{pmatrix} \in \begin{pmatrix} \partial f(x) + L^* u \\ -Lx \end{pmatrix} + \begin{pmatrix} 0 \\ \partial g^*(u) \end{pmatrix} + \begin{pmatrix} \nabla h(x) \\ 0 \end{pmatrix}$$



Davis-Yin splitting with
and swapping A and B

$$P = \begin{pmatrix} \frac{1}{\gamma} \text{Id} & -\tau L^* L & 0 \\ 0 & & \frac{1}{\tau} \text{Id} \end{pmatrix}$$

\equiv

Condat-Vu algorithm form II

$$\begin{cases} u^{k+1} = \text{prox}_{\tau g^*}(u^k + \tau Lx^k) \\ x^{k+1} = \text{prox}_{\gamma f}(x^k - \gamma \nabla h(x^k) - \gamma L^*(2u^{k+1} - u^k)) \end{cases}$$

Summary



nonsmooth functions



large scale



proximal splitting algorithms
designed using monotone and fixed-point
operator theory



Summary



nonsmooth functions



large scale



proximal splitting algorithms
designed using monotone and fixed-point
operator theory



LC et al. "Proximal Splitting Algorithms for Convex Optimization: A Tour of Recent Advances, with New Twists", *SIAM Review*, 2023

Salim, LC, Mishchenko, Richtárik, "Dualize, split, randomize: Toward Fast Nonsmooth Optimization Algorithms", *JOTA*, 2022

Convergence rates

Theorem – PD3O, with g continuous around Lx^* :

$$\psi(x^k) - \psi(x^*) = o(1/\sqrt{k})$$

Theorem – accelerated PD3O and PDDY
when h or f strongly convex, with varying stepsizes:

$$\|x^k - x^*\|^2 = O(1/k^2)$$

Theorem – linear convergence of PD3O and PDDY
when h or f strongly convex and g smooth:

$$\|x^k - x^*\|^2 \leq (1 - \rho)^k c_0$$

LC, Malinovsky, Richtárik, "Distributed Proximal Splitting Algorithms with Rates and Acceleration", 2022

PDDY algorithm

PDDY

input: initial points $x^0 \in \mathcal{X}$, $u^0 \in \mathcal{U}$;

stepsizes $\gamma > 0$, $\tau > 0$

$$h^0 := -L^* u^0$$

for $k = 0, 1, \dots$ **do**

$$\hat{x}^k := \text{prox}_{\gamma f} (x^k - \gamma \nabla h(x^k) + \gamma h^k)$$

$$u^{k+1} := \text{prox}_{\tau g^*} (u^k + \tau L \hat{x}^k)$$

$$h^{k+1} := -L^* u^{k+1}$$

$$x^{k+1} := \hat{x}^k + \gamma (h^{k+1} - h^k)$$

end for

PDDY algorithm

PDDY

input: initial points $x^0 \in \mathcal{X}$, $u^0 \in \mathcal{U}$;

stepsizes $\gamma > 0$, $\tau > 0$

$$h^0 := -L^* u^0$$

for $k = 0, 1, \dots$ **do**

$$\hat{x}^k := \text{prox}_{\gamma f} (x^k - \gamma \nabla h(x^k) + \gamma h^k)$$

$$u^{k+1} := \text{prox}_{\tau g^*} (u^k + \tau L \hat{x}^k)$$

$$h^{k+1} := -L^* u^{k+1}$$

$$x^{k+1} := \hat{x}^k + \gamma (h^{k+1} - h^k)$$

end for



$\text{prox}_{\tau g^*}$ can be costly

RandProx algorithm

RandProx

input: initial points $x^0 \in \mathcal{X}$, $u^0 \in \mathcal{U}$;

stepsizes $\gamma > 0$, $\tau > 0$; $\omega \geq 0$

$h^0 := -L^* u^0$

for $k = 0, 1, \dots$ **do**

$\hat{x}^k := \text{prox}_{\gamma f}(x^k - \gamma \nabla h(x^k) + \gamma h^k)$

$u^{k+1} := u^k + \frac{1}{1+\omega} \mathcal{R}^k(\text{prox}_{\tau g^*}(u^k + \tau L \hat{x}^k) - u^k)$

$h^{k+1} := -L^* u^{k+1}$

$x^{k+1} := \hat{x}^k + \gamma(1 + \omega)(h^{k+1} - h^k)$

end for

LC, Richtárik, "RandProx: Primal-Dual Optimization Algorithms with Randomized Proximal Updates", ICLR, 2023

RandProx algorithm

RandProx

input: initial points $x^0 \in \mathcal{X}$, $u^0 \in \mathcal{U}$;

stepsizes $\gamma > 0$, $\tau > 0$; $\omega \geq 0$

$h^0 := -L^* u^0$

for $k = 0, 1, \dots$ **do**

$\hat{x}^k := \text{prox}_{\gamma f}(x^k - \gamma \nabla h(x^k) + \gamma h^k)$

$u^{k+1} := u^k + \frac{1}{1+\omega} \mathcal{R}^k(\text{prox}_{\tau g^*}(u^k + \tau L \hat{x}^k) - u^k)$

$h^{k+1} := -L^* u^{k+1}$

$x^{k+1} := \hat{x}^k + \gamma(1 + \omega)(h^{k+1} - h^k)$

end for

$$\mathbb{E}[\mathcal{R}(r^k)] = r^k \quad \text{and} \quad \mathbb{E}[\|\mathcal{R}(r^k) - r^k\|^2] \leq \omega \|r^k\|^2$$

Linear convergence

Theorem 1. Suppose that $\mu_h > 0$ or $\mu_f > 0$, and $\mu_{g^*} > 0$.
For suitable γ and τ , $\forall k \geq 0$,

$$\mathbb{E}[\psi^k] \leq c^k \psi^0, \text{ where}$$

$$\psi^k = \frac{1}{\gamma} \left\| x^k - x^* \right\|^2 + (1 + \omega) \left(\frac{1}{\tau} + 2\mu_{g^*} \right) \left\| u^k - u^* \right\|^2,$$

$$c := \max \left(\frac{(1 - \gamma\mu_h)^2}{1 + \gamma\mu_f}, 1 - \frac{2\tau\mu_{g^*}}{(1 + \omega)(1 + 2\tau\mu_{g^*})} \right)$$

Examples

RandProx-skip

input: initial points $x^0 \in \mathcal{X}$, $u^0 \in \mathcal{U}$;

stepsizes $\gamma > 0$, $\tau > 0$; $p \in (0, 1]$

$h^0 := -L^* u^0$

for $k = 0, 1, \dots$ **do**

Flip a coin $\theta^k = (1 \text{ with probability } p, 0 \text{ else})$

if $\theta^k = 1$ **then**

$\hat{x}^k := \text{prox}_{\gamma f}(x^k - \gamma \nabla h(x^k) + \gamma h^k)$

$u^{k+1} := \text{prox}_{\tau g^*}(u^k + \tau L \hat{x}^k)$

$h^{k+1} := -L^* u^{k+1}$

$x^{k+1} := \hat{x}^k + \frac{\gamma}{p}(h^{k+1} - h^k)$

else

$x^{k+1} := \text{prox}_{\gamma f}(x^k - \gamma \nabla h(x^k) + \gamma h^k)$

$u^{k+1} := u^k$, $h^{k+1} := h^k$

end if

end for

$\mathcal{R}^t : r^t \mapsto$

$\begin{cases} \frac{1}{p} r^t & \text{with prob } p \\ 0 & \text{with prob } 1-p \end{cases}$

Examples

$$\min h + f + \sum_{m=1}^M g_m$$

RandProx-minibatch

input: initial points $x^0 \in \mathcal{X}$, $(u_m^0)_{m=1}^M \in \mathcal{X}^M$;

stepsize $\gamma > 0$; $s \in \{1, \dots, M\}$

$$h^0 := - \sum_{m=1}^M u_m^0$$

for $k = 0, 1, \dots$ **do**

$$\hat{x}^k := \text{prox}_{\gamma f}(x^k - \gamma \nabla h(x^k) + \gamma h^k)$$

pick $\Omega^k \subset \{1, \dots, M\}$ of size s unif. at random

for $m \in \Omega^k$ **do**

$$u_m^{k+1} := \text{prox}_{\frac{1}{\gamma M} g_m^*}(u_m^k + \frac{1}{\gamma M} \hat{x}^k)$$

end for

for $m \in \{1, \dots, M\} \setminus \Omega^k$ **do**

$$u_m^{k+1} := u_m^k$$

end for

$$h^{k+1} := - \sum_{m=1}^M u_m^{k+1}$$

$$x^{k+1} := \hat{x}^k + \frac{\gamma M}{s}(h^{k+1} - h^k)$$

end for

\mathcal{R}^k :
sampling

Examples

RandProx-FL

input: initial estimates $(x_i^0)_{i=1}^n \in \mathcal{X}^n$,
 $(u_i^0)_{i=1}^n \in \mathcal{X}^n$ such that $\sum_{i=1}^n u_i^0 = 0$;

stepsize $\gamma > 0$; $\omega \geq 0$

for $k = 0, 1, \dots$ **do**

for $i = 1, \dots, n$ at nodes in parallel **do**

$$\hat{x}_i^k := x_i^k - \gamma \nabla h_i(x_i^k) - \gamma u_i^k$$

$$a_i^k := \mathcal{R}^k(\hat{x}_i^k)$$

 // send compressed vector a_i^k to master

end for

$$a^k := \frac{1}{n} \sum_{i=1}^n a_i^k \quad // \text{aggregation at master}$$

 // broadcast a^k to all nodes

for $i = 1, \dots, n$ at nodes in parallel **do**

$$d_i^k := a_i^k - a^k$$

$$u_i^{k+1} := u_i^k + \frac{1}{\gamma(1+\omega)^2} d_i^k$$

$$x_i^{k+1} := \hat{x}_i^k - \frac{1}{1+\omega} d_i^k$$

end for

end for

$$\min \sum_{i=1}^n h_i$$

\mathcal{R}^k :

compression

Conclusion

nonsmooth large-scale optimization



proximal splitting algorithms



Perspectives

- ▶ acceleration
- ▶ different metric / Bregman distances
- ▶ cheaper random/inexact operations
- ▶ application to nonconvex optimization