

# RandProx: Primal–Dual Optimization Algorithms with Randomized Proximal Updates

[ICLR 2023]

Laurent Condat

Peter Richtárik

King Abdullah Univ. of  
Science and Technology  
(KAUST)  
Saudi Arabia



SIAM Conf. on Optim., Seattle, June 2, 2023

# Convex optimization

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \sum_{i=1}^n f_i(K_i x)$$

with

- linear operators  $K_i : \mathcal{X} \rightarrow \mathcal{U}_i$
- real Hilbert spaces  $\mathcal{X}, \mathcal{U}_i$
- convex functions  $f_i : \mathcal{U}_i \rightarrow \mathbb{R} \cup \{+\infty\}$



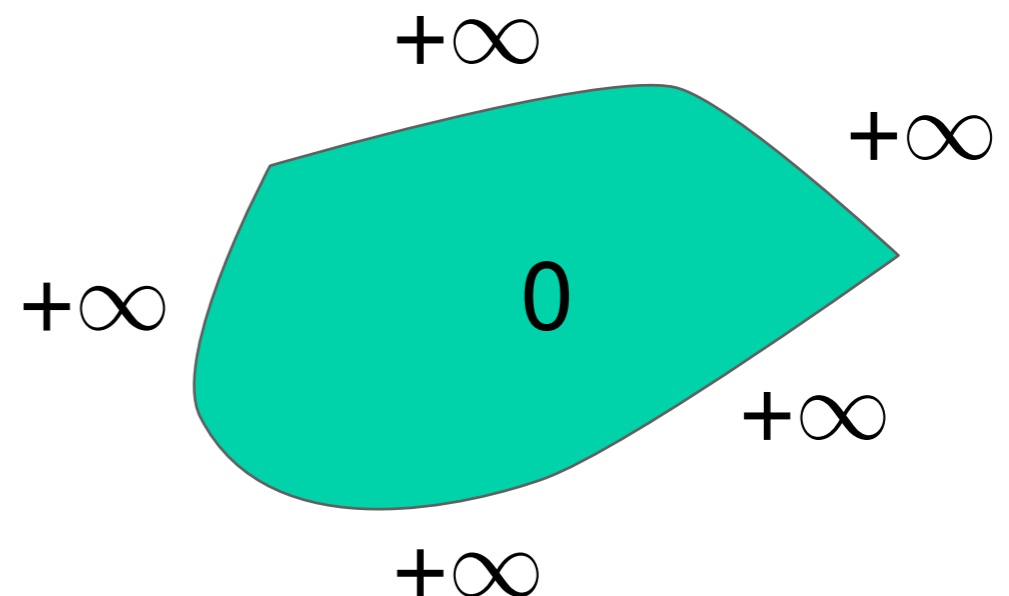
# Convex optimization

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \sum_{i=1}^n f_i(K_i x)$$

encompasses constraints:

$$\text{minimize}_{x \in \mathcal{X}} f(x) \quad \text{s.t. } x \in \Omega \quad \equiv \quad \text{minimize}_{x \in \mathcal{X}} f(x) + I_{\Omega}(x)$$

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

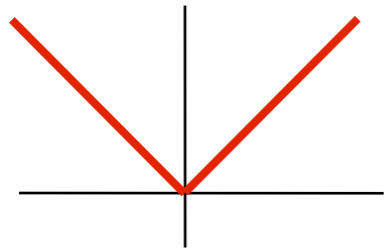


Note:  $I_{\Omega_1} + I_{\Omega_2} = I_{\Omega_1 \cap \Omega_2}$

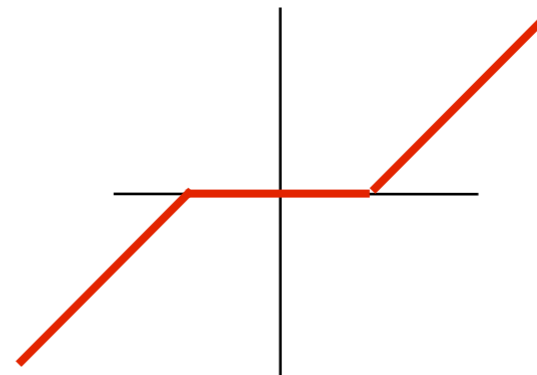


# The proximity operator

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



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



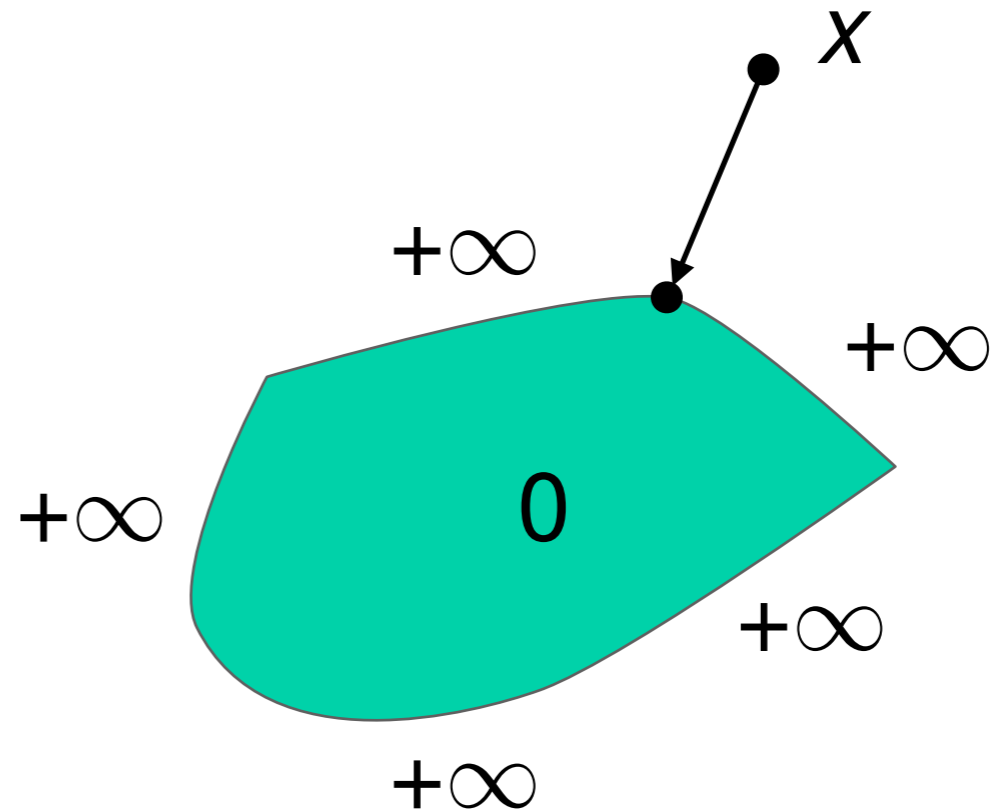
$$\text{prox}_f(x) = \text{sgn}(x) \max(|x| - 1, 0)$$



# The proximity operator

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

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



# The proximity operator

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

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

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

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

graph cuts, dynamic programming...

# Proximal splitting algorithms

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \sum_{i=1}^n f_i(K_i x)$$



No easy form of  $\text{prox}_{f_1+f_2}$  or  $\text{prox}_{f \circ K}$

# Proximal splitting algorithms

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \sum_{i=1}^n f_i(K_i x)$$



We want **full splitting**, with individual activation of  $K_i$ ,  $K_i^*$ , the gradient or proximity operator of  $f_i$ .





# Proximal splitting algorithms

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \left( f(x) + g(x) + \sum_{i=1}^n h_i(K_i x) \right)$$

with:

- $f$  smooth with  $L$ -Lipschitz grad  $\rightarrow$  calls to  $\nabla f$
- calls to  $\text{prox}_{\gamma g}$ ,  $\text{prox}_{\tau h_i}$ ,  $K_i$ ,  $K_i^*$



# Product space trick

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

$$h(\mathbf{u}) = \sum_{i=1}^n h_i(u_i)$$



$$h(\mathbf{K}x) = \sum_{i=1}^n h_i(K_i x)$$

$$\mathbf{K}x = (K_1 x, \dots, K_n x)$$

# Minimization of 3 functions

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \left( \underbrace{f(x)}_{\nabla f} + \underbrace{g(x)}_{\text{prox}_{\gamma} g} + \underbrace{h(Kx)}_{\text{prox}_{\tau} h}, K, K^* \right)$$

# Randomized algorithms

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

randomize  $\nabla f$

 SGD-type algorithms

# Randomized algorithms

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

↓

randomize  $\text{prox}_{\tau h}$

?

# The power of randomness

Find  $x^* \in \arg \min_{x \in \mathcal{X}} \sum_{i=1}^n f_i(x)$  using the  $\nabla f_i$

(every  $f_i$  is  $L$ -smooth and  $\mu$ -strongly convex)



lower bounds in Woodworth & Srebro [2016]

- deterministic algorithms:  $\Omega(n\sqrt{L/\mu} \log \epsilon^{-1})$
- randomized algorithms:  $\Omega((n + \sqrt{nL/\mu}) \log \epsilon^{-1})$

# Randomized algorithms

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

↓

randomize  $\text{prox}_{\tau h}$

?



# Proximal splitting algorithms

minimize  $f + g + h \circ K$

1979	$f + g$	👉	forward-backward alg.	
	$g + h$	👉	Douglas-Rachford alg. / ADMM	
2011	$g + h \circ K$	👉	Chambolle-Pock	
	$f + h \circ K$	👉	PAPC	
2013	$f + g + h \circ K$	👉	Condat, Vu	
2017	$f + g + h$	👉	Davis-Yin	
2018	$f + g + h \circ K$	👉	PD3O	Salim, LC et al., "Dualize, split, randomize: Fast nonsmooth optimization algorithms", <i>JOTA</i> , 2022
2020	$f + g + h \circ K$	👉	<b>PDDY</b> / AFBA	



# Proximal splitting algorithms

$$\text{minimize } f + g + h \circ K$$

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

Convergence speed:

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

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

$$\hat{x}^t := \text{prox}_{\gamma g} (x^t - \gamma \nabla f(x^t) - \gamma K^* u^t)$$

$$u^{t+1} := \text{prox}_{\tau h^*} (u^t + \tau K \hat{x}^t)$$

$$x^{t+1} := \hat{x}^t - \gamma K^* (u^{t+1} - u^t)$$

**end for**

---



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



# RandProx

---

## RandProx

---

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

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

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

$$\hat{x}^t := \text{prox}_{\gamma g} (x^t - \gamma \nabla f(x^t) - \gamma K^* u^t)$$

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

$$x^{t+1} := \hat{x}^t - \gamma(1 + \omega) K^* (u^{t+1} - u^t)$$

**end for**

---

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

$\mathcal{R}^t \equiv \text{Id}$ ,  $\omega = 0$   RandProx = PDDY

# Linear convergence

**Theorem 1.** *Suppose that  $\mu_f > 0$  or  $\mu_g > 0$ , and  $\mu_{h^*} > 0$ . For suitable  $\gamma$  and  $\tau$ ,  $\forall t \geq 0$ ,*

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

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

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

+ almost sure convergence



# Examples

## RandProx-skip

**input:** initial points  $x^0 \in \mathcal{X}$ ,  $u^0 \in \mathcal{U}$ ;  
stepsizes  $\gamma > 0$ ,  $\tau > 0$ ;  $p \in (0, 1]$

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

$$\hat{x}^t := \text{prox}_{\gamma g} (x^t - \gamma \nabla f(x^t) - \gamma K^* u^t)$$

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

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

$$u^{t+1} := \text{prox}_{\tau h^*} (u^t + \tau K \hat{x}^t)$$

$$x^{t+1} := \hat{x}^t - \frac{\gamma}{p} K^* (u^{t+1} - u^t)$$

**else**

$$u^{t+1} := u^t, x^{t+1} := \hat{x}^t$$

**end if**

**end for**

$$\mathcal{R}^t : d^t \mapsto \begin{cases} \frac{1}{p} d^t & \text{with prob } p \\ 0 & \text{with prob } 1-p \end{cases}$$

$$\omega = \frac{1}{p} - 1$$



# Examples

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

---

## RandProx-minibatch

---

**input:** initial points  $x^0 \in \mathcal{X}$ ,  $(u_i^0)_{i=1}^n \in \mathcal{X}^n$ ;

stepsize  $\gamma > 0$ ;  $k \in \{1, \dots, n\}$

$$v^0 := \sum_{i=1}^n u_i^0$$

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

$$\hat{x}^t := \text{prox}_{\gamma g}(x^t - \gamma \nabla f(x^t) - \gamma v^t)$$

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

**for**  $i \in \Omega^t$  **do**

$$u_i^{t+1} := \text{prox}_{\frac{1}{\gamma n} h_i^*}(u_i^t + \frac{1}{\gamma n} \hat{x}^t)$$

**end for**

**for**  $i \in \{1, \dots, n\} \setminus \Omega^t$  **do**

$$u_i^{t+1} := u_i^t$$

**end for**

$$v^{t+1} := \sum_{i=1}^n u_i^{t+1}$$

$$x^{t+1} := \hat{x}^t - \frac{\gamma n}{k}(v^{t+1} - v^t)$$

**end for**

---

$\mathcal{R}^t$ :  
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**  $t = 0, 1, \dots$  **do**

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

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

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

  // send compressed vector  $a_i^t$  to master

**end for**

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

  // broadcast  $a^t$  to all nodes

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

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

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

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

**end for**

**end for**

---

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

$\mathcal{R}^t$ :  
 compression



# Conclusion

A new **randomization technique** for PDDY,  
a generic primal-dual proximal splitting alg.





# Perspectives

► joint primal and dual randomization

LC et al. "TAMUNA: Doubly accelerated federated learning with local training, compression, and partial participation", preprint, 2023



# Perspectives

- ▶ joint primal and dual randomization

LC et al. "TAMUNA: Doubly accelerated federated learning with local training, compression, and partial participation", preprint, 2023

- ▶ different metric / Bregman distances

- ▶ acceleration?