

# MultiGrid Proximal Gradient Descent

**Andersen Ang**

ECS, USouthampton, UK

Homepage [angms.science](http://angms.science)

Andersen Ang, Hans De Sterck, Steve Vavasis,  
“MGProx: A nonsmooth multigrid proximal gradient method  
with adaptive restriction for strongly convex optimization”,  
SIAM Journal of Optimization, to appear, 2024

[arXiv 2302.04077](https://arxiv.org/abs/2302.04077) joint work with



Hans De Sterck



Steve Vavasis

25th International Symposium on Mathematical Programming, Montreal, Canada, July 26, 2024



## Setup

$$\operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{x})$$

- ▶ everything in finite dimensional Euclidean space

- ▶ not about

- ▶  $\mathbb{E}$
- ▶  $\nabla^2 f$
- ▶ linear

- ▶ big  $N$  in finite sum  $\frac{1}{N} \sum_{i=1}^N f_i(\mathbf{x})$

- ▶ application

- ▶  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is

- ▶  $\mu$ -strongly convex
- ▶  $L$ -smooth

- ▶  $g : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$  is

- ▶ convex
- ▶ proper
- ▶ possibly nonsmooth

- ▶ further assume a single point of non-differentiability

- ▶ separable:  $g(\mathbf{x}) = g_1(x_1) + g_2(x_2) + \dots$

- ▶ e.g.  $\max\{\cdot\}$ ,  $\|\cdot\|_1$

- ▶ proximable

What is the idea

► goal: solve

$$\operatorname{argmin}_{\mathbf{x}_{\text{Big}} \in \mathbb{R}^n} f_{\text{Big}}(\mathbf{x}_{\text{Big}}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}})$$

►  $n$  big, expensive

e.g. galaxy image (Lauga et. al., IML-FISTA.)



What is the idea

- ▶ goal: solve

$$\operatorname{argmin}_{\mathbf{x}_{\text{Big}} \in \mathbb{R}^n} f_{\text{Big}}(\mathbf{x}_{\text{Big}}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}})$$

- ▶  $n$  big, expensive

e.g. galaxy image (Lauga et. al., IML-FISTA.)



- ▶ natural idea: make use of *subspace*

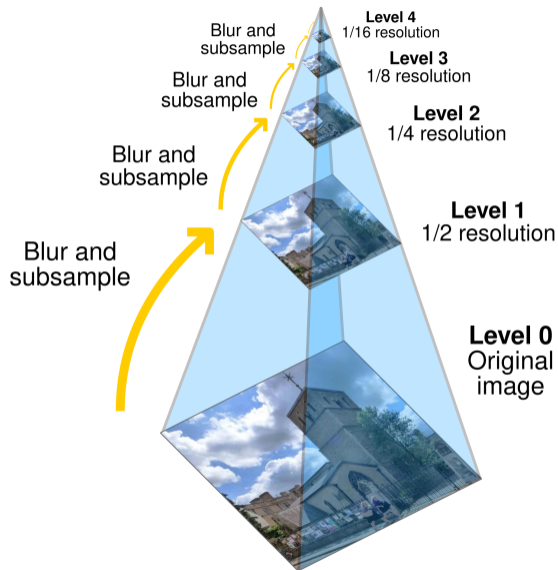
$$\operatorname{argmin}_{\mathbf{x}_{\text{small}}} f_{\text{small}}(\mathbf{x}_{\text{small}}) + g_{\text{small}}(\mathbf{x}_{\text{small}})$$

- ▶ how to define small problem?
- ▶ how to create small problem?
- ▶ how to solve small problem?

## Idea is old

- ▶ Caratheodory in 1910s
- ▶ Low rank by Eckart-Young in 1936
- ▶ PDEs: multigrid in 1962<sup>a</sup>
- ▶ CFD: model order reduction in 1967
- ▶ Linear programming: aggregation in 1977
- ▶ Computer vision: Pyramids in 1980s
- ▶ Natural language processing: Topic modelling in 1980s
- ▶ Wavelet in 1989
- ▶ Sparse linear statistic model in 1996
- ▶ Image segmentation: SuperPixel in 2000s
- ▶ Graph: compressing network by SuperNode in 2000s
- ▶ Core-set in 2010s
- ▶ Knowledge distillation

<sup>a</sup>R. P. Fedorenko, "A relaxation method for solving elliptic difference equations", USSR Computational Mathematics and Mathematical Physics, 1962



## What's there

- ▶ Multigrid for PDEs, a **whole field**
- ▶ mostly smooth problems
- ▶ two types
  - ▶ Full approximation scheme (FAS) /  $\tau$ -approximation scheme (1st-order method)
  - ▶ Newton-multigrid method (2nd-order method)
- ▶ works for nonsmooth problem exist, but
  - ▶ use smoothing, e.g.  $|x| \rightarrow \sqrt{x^2 + \epsilon}$
  - ▶ only simple constraint

## What's new

- ▶ Extend the FAS to nonsmooth
- ▶ No smoothing, face subdifferential directly
- ▶ Main output: proving it works
  - ▶ (new) adaptive restriction
  - ▶ consistent optimality between two worlds
  - ▶ descents
- ▶ other fancy stuffs
  - ▶ convergence rate  $1/k$
  - ▶ acceleration rate  $1/k^2$
  - ▶ good-looking curves

How to go from  $f_{\text{Big}}$  to  $f_{\text{small}}$ : Subspace

$$f_{\text{small}} \quad := \quad f_{\text{Big}} \circ \mathbf{R}$$

$$f_{\text{small}}(\mathbf{x}_{\text{small}}) \quad := \quad f_{\text{Big}}(\mathbf{R}\mathbf{x}_{\text{Big}})$$

$$\mathbf{x}_{\text{small}} \quad := \quad \mathbf{R}\mathbf{x}_{\text{Big}}$$

►  $\mathbf{R}$  is called *restriction* in PDEs

- $\mathbf{R}$  is a short-wide matrix that maps from  $\mathbb{R}^N$  to  $\mathbb{R}^n$
- if  $n < N$  short-wide: not one-to-one but many-to-one



How to go from  $f_{\text{Big}}$  to  $f_{\text{small}}$ : Subspace

$$\begin{aligned}f_{\text{small}} &:= f_{\text{Big}} \circ \mathbf{R} \\f_{\text{small}}(\mathbf{x}_{\text{small}}) &:= f_{\text{Big}}(\mathbf{R}\mathbf{x}_{\text{Big}}) \\ \mathbf{x}_{\text{small}} &:= \mathbf{R}\mathbf{x}_{\text{Big}}\end{aligned}$$

- ▶  $\mathbf{R}$  is called *restriction* in PDEs
  - ▶  $\mathbf{R}$  is a short-wide matrix that maps from  $\mathbb{R}^N$  to  $\mathbb{R}^n$
  - ▶ if  $n < N$  short-wide: not one-to-one but many-to-one
- ▶  $\mathbf{R}$  is a class of matrices
  - ▶ If  $\mathbf{R}$  is the  $i$ th row of  $\mathbf{I} \implies$  coordinate descent
  - ▶ If  $\mathbf{R}$  is random  $\implies$  random subspace method
  - ▶ If  $\mathbf{R}$  comes from PDE  $\implies$  multigrid method

How to go from  $f_{\text{Big}}$  to  $f_{\text{small}}$ : Subspace

$$\begin{aligned}f_{\text{small}} &:= f_{\text{Big}} \circ \mathbf{R} \\f_{\text{small}}(\mathbf{x}_{\text{small}}) &:= f_{\text{Big}}(\mathbf{R}\mathbf{x}_{\text{Big}}) \\ \mathbf{x}_{\text{small}} &:= \mathbf{R}\mathbf{x}_{\text{Big}}\end{aligned}$$

- ▶  $\mathbf{R}$  is called *restriction* in PDEs
  - ▶  $\mathbf{R}$  is a short-wide matrix that maps from  $\mathbb{R}^N$  to  $\mathbb{R}^n$
  - ▶ if  $n < N$  short-wide: not one-to-one but many-to-one
- ▶  $\mathbf{R}$  is a class of matrices
  - ▶ If  $\mathbf{R}$  is the  $i$ th row of  $\mathbf{I} \implies$  coordinate descent
  - ▶ If  $\mathbf{R}$  is random  $\implies$  random subspace method
  - ▶ If  $\mathbf{R}$  comes from PDE  $\implies$  multigrid method
- ▶ No-free-lunch
  - ▶ How to choose / design  $\mathbf{R} \implies$  by experience, open problem
  - ▶ I am “cheating” by using a “known”  $\mathbf{R}$

FAS multigrid (red are new) for  $\operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) + g(\mathbf{x})$

1.  $\mathbf{x}_{\text{Big}}^{k+1/2} = \operatorname{prox}_{\alpha g}(\mathbf{x}_{\text{Big}}^k - \alpha \nabla f(\mathbf{x}_{\text{Big}}^k))$

ProxGD in Big world

FAS multigrid (red are new) for  $\operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) + g(\mathbf{x})$

1.  $\mathbf{x}_{\text{Big}}^{k+1/2} = \operatorname{prox}_{\alpha g} \left( \mathbf{x}_{\text{Big}}^k - \alpha \nabla f(\mathbf{x}_{\text{Big}}^k) \right)$

ProxGD in Big world

2.  $\mathbf{x}_{\text{small}}^{k+1/2} = \mathbf{R} \mathbf{x}_{\text{Big}}^{k+1/2}$

Restriction: Big-to-small

$$f_{\text{small}} = f_{\text{Big}} \circ \mathbf{R}, \quad g_{\text{small}} = g_{\text{Big}} \circ \mathbf{R}$$

adaptive  $\mathbf{R}$  is new

FAS multigrid (red are new) for  $\operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) + g(\mathbf{x})$

$$1. \quad \mathbf{x}_{\text{Big}}^{k+1/2} = \operatorname{prox}_{\alpha g} \left( \mathbf{x}_{\text{Big}}^k - \alpha \nabla f(\mathbf{x}_{\text{Big}}^k) \right)$$

ProxGD in Big world

$$2. \quad \mathbf{x}_{\text{small}}^{k+1/2} = \mathbf{R} \mathbf{x}_{\text{Big}}^{k+1/2}$$

Restriction: Big-to-small

$$f_{\text{small}} = f_{\text{Big}} \circ \mathbf{R}, \quad g_{\text{small}} = g_{\text{Big}} \circ \mathbf{R}$$

adaptive  $\mathbf{R}$  is new

$$3. \quad \boldsymbol{\tau} \in \partial \left( f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) + g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \right) \ominus \mathbf{R} \partial \left( f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) \right)$$

WTF

inclusion & Minkowski sum

FAS multigrid (red are new) for  $\operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) + g(\mathbf{x})$

$$1. \mathbf{x}_{\text{Big}}^{k+1/2} = \operatorname{prox}_{\alpha g} \left( \mathbf{x}_{\text{Big}}^k - \alpha \nabla f(\mathbf{x}_{\text{Big}}^k) \right)$$

ProxGD in Big world

$$2. \mathbf{x}_{\text{small}}^{k+1/2} = \mathbf{R} \mathbf{x}_{\text{Big}}^{k+1/2}$$

Restriction: Big-to-small

$$f_{\text{small}} = f_{\text{Big}} \circ \mathbf{R}, \quad g_{\text{small}} = g_{\text{Big}} \circ \mathbf{R}$$

adaptive  $\mathbf{R}$  is new

$$3. \boldsymbol{\tau} \in \partial \left( f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) + g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \right) \ominus \mathbf{R} \partial \left( f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) \right)$$

WTF

inclusion & Minkowski sum

$$4. \mathbf{x}_{\text{small}}^{k+1} = \operatorname{prox}_{\alpha g} \left( \mathbf{x}_{\text{small}}^{k+1/2} - \alpha (\nabla f(\mathbf{x}_{\text{small}}^{k+1/2}) - \boldsymbol{\tau}) \right)$$

ProxGD in small world

FAS multigrid (red are new) for  $\operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) + g(\mathbf{x})$

$$1. \quad \mathbf{x}_{\text{Big}}^{k+1/2} = \operatorname{prox}_{\alpha g} \left( \mathbf{x}_{\text{Big}}^k - \alpha \nabla f(\mathbf{x}_{\text{Big}}^k) \right)$$

ProxGD in Big world

$$2. \quad \mathbf{x}_{\text{small}}^{k+1/2} = \mathbf{R} \mathbf{x}_{\text{Big}}^{k+1/2}$$

Restriction: Big-to-small

$$f_{\text{small}} = f_{\text{Big}} \circ \mathbf{R}, \quad g_{\text{small}} = g_{\text{Big}} \circ \mathbf{R}$$

adaptive  $\mathbf{R}$  is new

$$3. \quad \boldsymbol{\tau} \in \partial \left( f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) + g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \right) \ominus \mathbf{R} \partial \left( f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) \right)$$

WTF

inclusion & Minkowski sum

$$4. \quad \mathbf{x}_{\text{small}}^{k+1} = \operatorname{prox}_{\alpha g} \left( \mathbf{x}_{\text{small}}^{k+1/2} - \alpha (\nabla f(\mathbf{x}_{\text{small}}^{k+1/2}) - \boldsymbol{\tau}) \right)$$

ProxGD in small world

$$5. \quad \mathbf{x}_{\text{Big}}^{k+1} = \mathbf{x}_{\text{Big}}^{k+1/2} + \alpha \mathbf{P} \left( \mathbf{x}_{\text{small}}^{k+1} - \mathbf{x}_{\text{small}}^{k+1/2} \right)$$

use small to correct Big

no prox

existence of  $\alpha$  is new

## Where are the fun things

- ▶ Subdifferential can be  $\emptyset$ , crazy things will happen

We fixed this: empty-set will not appear

- ▶ How to choose  $\tau$ ?

We fixed this: the algo always works regardless of choice of  $\tau$

- ▶  $P(\mathbf{x}_{\text{small}}^{k+1} - \mathbf{x}_{\text{small}}^{k+1/2})$  works?

We fixed this: it always is a descent direction for  $x_{\text{Big}}$

- ▶  $\alpha$  exists?

We fixed this:  $\alpha > 0$  exists



So ... what is new?

$$\boldsymbol{\tau} \in \partial\left(f_{\text{small}}(\boldsymbol{x}_{\text{small}}^{k+1/2}) + g_{\text{small}}(\boldsymbol{x}_{\text{small}}^{k+1/2})\right) \ominus \mathbf{R}\partial\left(f_{\text{Big}}(\boldsymbol{x}_{\text{Big}}^{k+1/2}) + g_{\text{Big}}(\boldsymbol{x}_{\text{Big}}^{k+1/2})\right)$$

Minkowski sum

So ... what is new?

$$\begin{aligned}\tau &\in \partial\left(f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) + g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2})\right) \ominus \mathbf{R}\partial\left(f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})\right) \\ &= \partial f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \oplus \partial g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \ominus \mathbf{R}\left(\partial f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) \oplus \partial g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})\right)\end{aligned}$$

Minkowski sum

Moreau-Rockafellar thm. (new<sup>1</sup>)

So ... what is new?

$$\begin{aligned}\tau &\in \partial\left(f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) + g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2})\right) \ominus \mathbf{R}\partial\left(f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})\right) \\ &= \partial f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \oplus \partial g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \ominus \mathbf{R}\left(\partial f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) \oplus \partial g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})\right) \\ &= \nabla f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) + \partial g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \ominus \mathbf{R}\left(\nabla f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) + \partial g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})\right)\end{aligned}$$

Minkowski sum

Moreau-Rockafellar thm. (new<sup>1</sup>)

diff.able part becomes singleton

So ... what is new?

$$\begin{aligned}
 \tau &\in \partial\left(f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) + g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2})\right) \ominus \mathbf{R}\partial\left(f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})\right) \\
 &= \partial f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \oplus \partial g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \ominus \mathbf{R}\left(\partial f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) \oplus \partial g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})\right) \\
 &= \nabla f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) + \partial g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \ominus \mathbf{R}\left(\nabla f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) + \partial g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})\right) \\
 &= \underbrace{\left\{\nabla f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) - \mathbf{R}\nabla f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})\right\}}_{\text{singleton, simple part, exists in literature}} \oplus \underbrace{\left(\partial g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \oplus (-\mathbf{R})\partial g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})\right)}_{\text{set Minkowski sum, hard part, new thing}}
 \end{aligned}$$

Minkowski sum

Moreau-Rockafellar thm. (new<sup>1</sup>)

diff.able part becomes singleton

So ... what is new?

$$\begin{aligned}
 \tau &\in \partial\left(f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) + g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2})\right) \ominus \mathbf{R}\partial\left(f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})\right) && \text{Minkowski sum} \\
 &= \partial f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \oplus \partial g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \ominus \mathbf{R}\left(\partial f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) \oplus \partial g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})\right) && \text{Moreau-Rockafellar thm. (new}^1) \\
 &= \nabla f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) + \partial g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \ominus \mathbf{R}\left(\nabla f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) + \partial g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})\right) && \text{diff.able part becomes singleton} \\
 &= \underbrace{\left\{\nabla f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) - \mathbf{R}\nabla f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})\right\}}_{\text{singleton, simple part, exists in literature}} \oplus \underbrace{\left\{\partial g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \oplus (-\mathbf{R})\partial g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})\right\}}_{\text{set Minkowski sum, hard part, new thing}} \\
 &= \underbrace{\left\{\nabla f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) - \mathbf{R}\nabla f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})\right\}}_{\text{singleton, simple part, exists in literature}} \oplus \partial g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) + (-\mathbf{R})\partial g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) && \text{adaptive restriction (new}^2)
 \end{aligned}$$

1. We made this true
2. To make my life easier:  $\mathbf{R}$  maps all set-valued entries of  $\partial g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})$  to the singleton  $\{0\}$   
**Open problem:** keep dealing with  $\oplus(-\mathbf{R})$  instead of  $\oplus(-\mathbf{R})$
3. the whole scheme is different from existing literature

What's the heck is  $\tau$ ?

►  $\tau$  links Big-world and small-world

What's the heck is  $\tau$ ?

►  $\tau$  links Big-world and small-world

► **Theorem** IF  $\mathbf{x}_{\text{Big}}$  solves  $\operatorname{argmin}_{\mathbf{x}_{\text{Big}} \in \mathbb{R}^n} f_{\text{Big}}(\mathbf{x}_{\text{Big}}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}})$

THEN  $R\mathbf{x}_{\text{Big}}$  solves  $\operatorname{argmin}_{\mathbf{x}_{\text{small}} \in \mathbb{R}^n} f_{\text{small}}(\mathbf{x}_{\text{small}}) + g_{\text{small}}(\mathbf{x}_{\text{small}}) - \langle \boldsymbol{\tau}, \mathbf{x}_{\text{small}} \rangle$

What's the heck is  $\tau$ ?

►  $\tau$  links Big-world and small-world

► **Theorem** IF  $\mathbf{x}_{\text{Big}}$  solves  $\operatorname{argmin}_{\mathbf{x}_{\text{Big}} \in \mathbb{R}^n} f_{\text{Big}}(\mathbf{x}_{\text{Big}}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}})$

THEN  $R\mathbf{x}_{\text{Big}}$  solves  $\operatorname{argmin}_{\mathbf{x}_{\text{small}} \in \mathbb{R}^n} f_{\text{small}}(\mathbf{x}_{\text{small}}) + g_{\text{small}}(\mathbf{x}_{\text{small}}) - \langle \boldsymbol{\tau}, \mathbf{x}_{\text{small}} \rangle$

(In algo) ProxGrad converges in Big-world  $\iff$  ProxGrad converges in small-world



What's the heck is  $\tau$ ?

►  $\tau$  links Big-world and small-world

► **Theorem** IF  $\mathbf{x}_{\text{Big}}$  solves  $\operatorname{argmin}_{\mathbf{x}_{\text{Big}} \in \mathbb{R}^n} f_{\text{Big}}(\mathbf{x}_{\text{Big}}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}})$

THEN  $R\mathbf{x}_{\text{Big}}$  solves  $\operatorname{argmin}_{\mathbf{x}_{\text{small}} \in \mathbb{R}^n} f_{\text{small}}(\mathbf{x}_{\text{small}}) + g_{\text{small}}(\mathbf{x}_{\text{small}}) - \langle \tau, \mathbf{x}_{\text{small}} \rangle$

(In algo) ProxGrad converges in Big-world  $\iff$  ProxGrad converges in small-world

► I am NOT saying:  $\mathbf{x}_{\text{Big}}$  solves Big-problem  $\iff$   $\mathbf{x}_{\text{small}}$  solves small-problem

► I am saying:  $\mathbf{x}_{\text{Big}}$  solves Big-problem  $\iff$   $\mathbf{x}_{\text{small}}$  solves small-perturbed-problem

►  $\tau$  is Big-perturb-small so that {sol of small-problem} flavours {sol of Big-problem}

What's the heck is  $\tau$ ?

►  $\tau$  links Big-world and small-world

► **Theorem** IF  $\mathbf{x}_{\text{Big}}$  solves  $\operatorname{argmin}_{\mathbf{x}_{\text{Big}} \in \mathbb{R}^n} f_{\text{Big}}(\mathbf{x}_{\text{Big}}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}})$

THEN  $\mathbf{R}\mathbf{x}_{\text{Big}}$  solves  $\operatorname{argmin}_{\mathbf{x}_{\text{small}} \in \mathbb{R}^n} f_{\text{small}}(\mathbf{x}_{\text{small}}) + g_{\text{small}}(\mathbf{x}_{\text{small}}) - \langle \boldsymbol{\tau}, \mathbf{x}_{\text{small}} \rangle$

(In algo) ProxGrad converges in Big-world  $\iff$  ProxGrad converges in small-world

► I am NOT saying:  $\mathbf{x}_{\text{Big}}$  solves Big-problem  $\iff$   $\mathbf{x}_{\text{small}}$  solves small-problem

► I am saying:  $\mathbf{x}_{\text{Big}}$  solves Big-problem  $\iff$   $\mathbf{x}_{\text{small}}$  solves small-perturbed-problem

►  $\tau$  is Big-perturb-small so that {sol of small-problem} flavours {sol of Big-problem}

► How to prove: definition of  $\tau$ ,  $\mathbf{R}$ , convexity of obj functions, 1st-order subdiff. optimality

Coarse correction step: inequality

$$\mathbf{x}_{\text{Big}}^{k+1} = \mathbf{x}_{\text{Big}}^{k+1/2} + \alpha \mathbf{P}(\mathbf{x}_{\text{small}}^{k+1} - \mathbf{x}_{\text{small}}^{k+1/2})$$

► **Theorem**  $\mathbf{P}(\mathbf{x}_{\text{small}}^{k+1} - \mathbf{x}_{\text{small}}^{k+1/2})$  is a descent direction in Big-world

$$\left\langle \partial \left( f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) \right), \mathbf{P}(\mathbf{x}_{\text{small}}^{k+1}(\boldsymbol{\tau}) - \mathbf{x}_{\text{small}}^{k+1/2}) \right\rangle < 0$$

the inequality is strict

Coarse correction step: inequality

$$\mathbf{x}_{\text{Big}}^{k+1} = \mathbf{x}_{\text{Big}}^{k+1/2} + \alpha \mathbf{P}(\mathbf{x}_{\text{small}}^{k+1} - \mathbf{x}_{\text{small}}^{k+1/2})$$

► **Theorem**  $\mathbf{P}(\mathbf{x}_{\text{small}}^{k+1} - \mathbf{x}_{\text{small}}^{k+1/2})$  is a descent direction in Big-world

$$\left\langle \partial \left( f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) \right), \mathbf{P}(\mathbf{x}_{\text{small}}^{k+1}(\tau) - \mathbf{x}_{\text{small}}^{k+1/2}) \right\rangle < 0$$

the inequality is strict

► The inequality holds for

► **any**  $\tau$  you choose to get  $\mathbf{x}_{\text{small}}^{k+1}(\tau)$

► **any** subgradient in  $\partial \left( f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) \right)$

Coarse correction step: inequality

$$\mathbf{x}_{\text{Big}}^{k+1} = \mathbf{x}_{\text{Big}}^{k+1/2} + \alpha \mathbf{P}(\mathbf{x}_{\text{small}}^{k+1} - \mathbf{x}_{\text{small}}^{k+1/2})$$

► **Theorem**  $\mathbf{P}(\mathbf{x}_{\text{small}}^{k+1} - \mathbf{x}_{\text{small}}^{k+1/2})$  is a descent direction in Big-world

$$\left\langle \partial \left( f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) \right), \mathbf{P}(\mathbf{x}_{\text{small}}^{k+1}(\boldsymbol{\tau}) - \mathbf{x}_{\text{small}}^{k+1/2}) \right\rangle < 0$$

the inequality is strict

► The inequality holds for

► **any**  $\boldsymbol{\tau}$  you choose to get  $\mathbf{x}_{\text{small}}^{k+1}(\boldsymbol{\tau})$

► **any** subgradient in  $\partial \left( f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) \right)$

► How to prove: definition & convexity

Coarse correction step: stepsize exists

$$\mathbf{x}_{\text{Big}}^{k+1} = \mathbf{x}_{\text{Big}}^{k+1/2} + \alpha \mathbf{P}(\mathbf{x}_{\text{small}}^{k+1} - \mathbf{x}_{\text{small}}^{k+1/2}) \quad (\text{A})$$

► **Theorem**  $\alpha > 0$  exists such that  $F_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1}) < F_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})$  ( $F := f + g$ )

► **3-sentence proof**

1. The strict inequality in

$$\left\langle \partial F_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}), \mathbf{P}(\mathbf{x}_{\text{small}}^{k+1}(\boldsymbol{\tau}) - \mathbf{x}_{\text{small}}^{k+1/2}) \right\rangle < 0$$

means  $\partial F_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})$  is strictly inside a half-space with normal  $\mathbf{N} := \mathbf{P}(\mathbf{x}_{\text{small}}^{k+1}(\boldsymbol{\tau}) - \mathbf{x}_{\text{small}}^{k+1/2})$

2. Subdifferential is a compact convex set  $\xRightarrow{1}$  strict separation  $\implies \partial F_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})$  must be a positive distance (a fact  $\alpha > 0$ ) from that hyperplane defined by  $\mathbf{N}$ .

3. Evaluate the support func. of  $\partial F_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2})$ , i.e., the directional derivative of  $F_{\text{Big}}$  at  $\mathbf{x}_{\text{Big}}^{k+1/2}$  in the direction  $\mathbf{N}$ , we are done

► Sad news: we only have descent condition, not sufficient descent condition

► Deep thing in the compactness of subdifferential

## Theoretical results

$$1. \mathbf{x}_{\text{Big}}^{k+1/2} = \text{prox}_{\alpha g} \left( \mathbf{x}_{\text{Big}}^k - \alpha \nabla f(\mathbf{x}_{\text{Big}}^k) \right)$$

$$2. \mathbf{x}_{\text{small}}^{k+1/2} = \mathbf{R} \mathbf{x}_{\text{Big}}^{k+1/2}$$

$$3. \boldsymbol{\tau} \in \partial \left( f_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) + g_{\text{small}}(\mathbf{x}_{\text{small}}^{k+1/2}) \right) \\ \ominus \mathbf{R} \partial \left( f_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) + g_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}) \right)$$

$$4. \mathbf{x}_{\text{small}}^{k+1} = \text{prox}_{\alpha g} \left( \mathbf{x}_{\text{small}}^{k+1/2} - \alpha (\nabla f(\mathbf{x}_{\text{small}}^{k+1/2}) - \boldsymbol{\tau}) \right)$$

$$5. \mathbf{x}_{\text{Big}}^{k+1} = \mathbf{x}_{\text{Big}}^{k+1/2} + \alpha \mathbf{P} \left( \mathbf{x}_{\text{small}}^{k+1} - \mathbf{x}_{\text{small}}^{k+1/2} \right)$$

1. At convergence,  $\mathbf{x}_{\ell}^k$  has a fixed-pt. property  $\forall$  level  $\ell$

2. **Nonsmooth** angle condition

$$\left\langle \partial F_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}), \mathbf{P}(\mathbf{x}_{\text{small}}^{k+1}(\boldsymbol{\tau}) - \mathbf{x}_{\text{small}}^{k+1/2}) \right\rangle < 0$$

3. Descent property: stepsize  $\alpha > 0$  exists and

$\mathbf{P}(\mathbf{x}_{\text{small}}^{k+1}(\boldsymbol{\tau}) - \mathbf{x}_{\text{small}}^{k+1/2})$  is a descent direction at  $\mathbf{x}_{\text{Big}}^{k+1/2}$

$$F_{\text{Big}} \left( \mathbf{x}_{\text{Big}}^{k+1/2} + \alpha \mathbf{P}(\mathbf{x}_{\text{small}}^{k+1}(\boldsymbol{\tau}) - \mathbf{x}_{\text{small}}^{k+1/2}) \right) < F_{\text{Big}}(\mathbf{x}_{\text{Big}}^{k+1/2}).$$

4.  $\{F_0(\mathbf{x}_{\text{Big}}^k)\}_{k \in \mathbb{N}}$  converges to  $F_{\text{Big}}^* := \inf F_{\text{Big}}$ , with

► a sublinear rate

$$F_{\text{Big}}(\mathbf{x}_{\text{Big}}^k) - F_{\text{Big}}^* \leq \frac{c}{k}$$

► a linear rate

$$F_{\text{Big}}(\mathbf{x}_{\text{Big}}^k) - F_{\text{Big}}^* \leq c \left( 1 - \frac{\mu_{\text{Big}}}{L_{\text{Big}}} \right)^k.$$

► with acceleration

$$F_{\text{Big}}(\mathbf{x}_{\text{Big}}^k) - F_{\text{Big}}^* \leq \frac{d}{ak^2 + bk + c}$$

5.  $\{\mathbf{x}_{\text{Big}}^k\}_{k \in \mathbb{N}} \xrightarrow{k} \mathbf{x}_{\text{Big}}^*$

# Why stop at 2-level?

**Algorithm 4.1**  $L$ -level MGProx with V-cycle structure for an approximate solution of (1.1)

Initialize  $x_0^1$  and the full version of  $R_{\ell \rightarrow \ell+1}, P_{\ell+1 \rightarrow \ell}$  for  $\ell \in \{0, 1, \dots, L-1\}$

**for**  $k = 1, 2, \dots$  **do**

Set  $\tau_{-1 \rightarrow 0}^{k+1} = 0$

**for**  $\ell = 0, 1, \dots, L-1$  **do**

$$y_\ell^{k+1} = \text{prox}_{\frac{1}{L_\ell} g_\ell} \left( x_\ell^k - \frac{\nabla f_\ell(x_\ell^k) - \tau_{\ell-1 \rightarrow \ell}^{k+1}}{L_\ell} \right) \quad \text{pre-smoothing}$$

$$x_{\ell+1}^k = R_{\ell \rightarrow \ell+1}(y_\ell^{k+1}) y_\ell^{k+1} \quad \text{restriction to next level}$$

$$\tau_{\ell \rightarrow \ell+1}^{k+1} \in \partial F_{\ell+1}(x_{\ell+1}^k) - R_{\ell \rightarrow \ell+1}(y_\ell^{k+1}) \partial F_\ell(y_\ell^{k+1}) \quad \text{create tau vector}$$

**end for**

$$w_L^{k+1} = \underset{\xi}{\text{argmin}} \left\{ F_L^\tau(\xi) := F_L(\xi) - \langle \tau_{L-1 \rightarrow L}^{k+1}, \xi \rangle \right\} \quad \text{solve the level-}L \text{ coarse problem}$$

**for**  $\ell = L-1, L-2, \dots, 0$  **do**

$$z_\ell^{k+1} = y_\ell^{k+1} + \alpha P_{\ell+1 \rightarrow \ell}(w_{\ell+1}^{k+1} - x_{\ell+1}^k) \quad \text{coarse correction}$$

$$w_\ell^{k+1} = \text{prox}_{\frac{1}{L_\ell} g_\ell} \left( z_\ell^{k+1} - \frac{\nabla f_\ell(z_\ell^{k+1}) - \tau_{\ell-1 \rightarrow \ell}^{k+1}}{L_\ell} \right) \quad \text{post-smoothing}$$

**end for**

$$x_0^{k+1} = w_0^{k+1} \quad \text{update the fine variable}$$

**end for**

$\ell$  level

$k$  iteration counter

$x_\ell$  variable at level  $\ell$  before ProxGrad

$y_\ell$  variable at level  $\ell$  after ProxGrad

$f_\ell$  smooth part at level  $\ell$

$g_\ell$  nonsmooth part at level  $\ell$

$L_\ell$  Lipschitz constant of  $\nabla f_\ell$  at level  $\ell$

$L$  number of levels

► Reduction in problem size  $n_0 \rightarrow \frac{1}{4}n_0 \rightarrow \frac{1}{16}n_0 \rightarrow \frac{1}{64}n_0 \rightarrow \frac{1}{256}n_0 \rightarrow \frac{1}{1024}n_0$

► Per-iteration cost by geometric series  $a + ar + ar^2 + \dots \rightarrow \frac{a}{1-r}$ . For  $r = \frac{1}{4}$ , V-cycle is  $2.66n_0$  for all single proximal gradient update.



## Experiment

- ▶ Tons of papers on multigrid methods for applications

## Experiment

- ▶ Tons of papers on multigrid methods for applications
- ▶ Elastic Obstacle Problem:

$$\min_u \iint_{\Omega} \sqrt{1 + \|\nabla u\|_{L^2}^2} dx dy + \lambda \iint_{\Omega} \|(\phi - u)_+\|_{L^1} dx dy \quad \text{s.t.} \quad u = 0 \text{ on } \partial\Omega,$$

$$\operatorname{argmin}_{\mathbf{u} \in \mathbb{R}^{N^2}} h^2 \sum_{i=1}^N \sum_{j=1}^N \sqrt{1 + (\mathbf{D}_{(i,j),:} \mathbf{u})^2 + (\mathbf{E}_{(i,j),:} \mathbf{u})^2} + h^2 \lambda \|(\phi - \mathbf{u})_+\|_1$$

$$\operatorname{argmin}_{\mathbf{u} \in \mathbb{R}^{N^2}} f_0(\mathbf{u}) + g_0(\mathbf{u}) := \sum_{i=1}^N \sum_{j=1}^N \psi(\mathbf{F}_{(i,j),:} \mathbf{u}) + \lambda \|(\phi - \mathbf{u})_+\|_1, \quad (\text{EOP})$$

$$\psi : \mathbb{R}^2 \rightarrow \mathbb{R} : (s, t) \mapsto \sqrt{1 + s^2 + t^2}, \quad \mathbf{F}_{(i,j),:} := \begin{bmatrix} \mathbf{D}_{(i,j),:} \\ \mathbf{E}_{(i,j),:} \end{bmatrix}.$$

## Experiment

- ▶ Tons of papers on multigrid methods for applications
- ▶ Elastic Obstacle Problem:

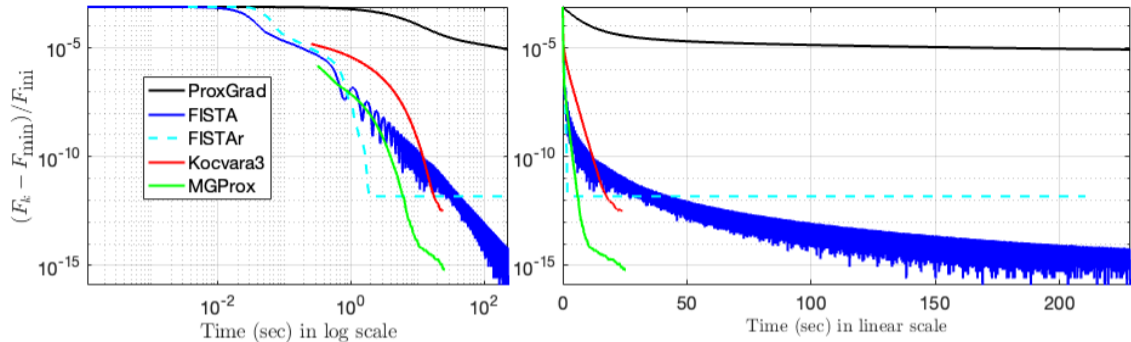
$$\min_u \iint_{\Omega} \sqrt{1 + \|\nabla u\|_{L^2}^2} dx dy + \lambda \iint_{\Omega} \|(\phi - u)_+\|_{L^1} dx dy \quad \text{s.t.} \quad u = 0 \text{ on } \partial\Omega,$$

$$\operatorname{argmin}_{\mathbf{u} \in \mathbb{R}^{N^2}} h^2 \sum_{i=1}^N \sum_{j=1}^N \sqrt{1 + (\mathbf{D}_{(i,j),:} \mathbf{u})^2 + (\mathbf{E}_{(i,j),:} \mathbf{u})^2} + h^2 \lambda \|(\phi - \mathbf{u})_+\|_1$$

$$\operatorname{argmin}_{\mathbf{u} \in \mathbb{R}^{N^2}} f_0(\mathbf{u}) + g_0(\mathbf{u}) := \sum_{i=1}^N \sum_{j=1}^N \psi(\mathbf{F}_{(i,j),:} \mathbf{u}) + \lambda \|(\phi - \mathbf{u})_+\|_1, \quad (\text{EOP})$$

$$\psi : \mathbb{R}^2 \rightarrow \mathbb{R} : (s, t) \mapsto \sqrt{1 + s^2 + t^2}, \quad \mathbf{F}_{(i,j),:} := \begin{bmatrix} \mathbf{D}_{(i,j),:} \\ \mathbf{E}_{(i,j),:} \end{bmatrix}.$$

- ▶ **Theorem** EOP problem is  $\mu$ -strongly convex and smooth part is  $L$ -smooth
  - ▶ global tight  $L$  is unknown (to us)
  - ▶  $\mu > 0$  is unknown
  - ▶  $\mu \rightarrow 0$



$(\mu, L)$  unknown  $\implies$  cannot do parameter-based restart

For  $(2^8 - 1)^2 = 65025$  number of variable

Method	iterations $k$	time (sec.)	$(F_0(x_0^k) - F_0^{\min})/F_0(x_0^{\text{ini}})$
ProxGrad	$> 10^5$	335.9	$8.33 \times 10^{-8}$
FISTA	$> 10^5$	332.62	$6.64 \times 10^{-8}$
FISTA-r	$> 10^5$	364.89	$6.64 \times 10^{-8}$
Kocvara3 $N_s = 100$	$> 10^3$	986.53	$8.07 \times 10^{-8}$
MGProx $N_s = 100$	50	48.37	$1.32 \times 10^{-10}$

## Last page - summary

Andersen Ang, Hans De Sterck, Steve Vavasis,  
“MGProx: A nonsmooth multigrid proximal gradient method with adaptive restriction for  
strongly convex optimization”,  
SIAM Journal of Optimization, to appear, 2024  
[arXiv 2302.04077](https://arxiv.org/abs/2302.04077)

- ▶ discussed some funny things
- ▶ many open problems / extensions / improvement opportunities
- ▶ Ads: I actually have PhD positions

End of document