

Newton methods for structured nonsmooth optimization

proximal identification and fast local convergence

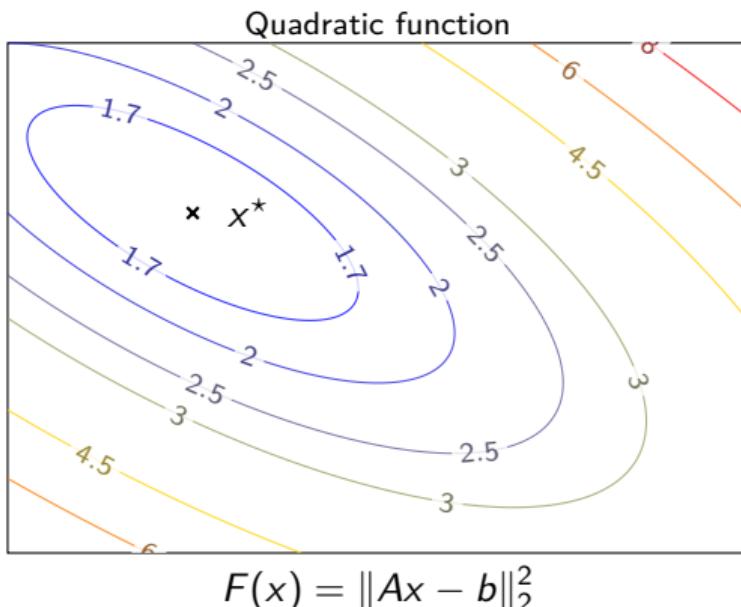
Gilles Bareilles
LJK, Univ. Grenoble Alpes
gbareilles.fr
December 2, 2022

Chair	Nadia Brauner
Reviewer	Jalal Fadili
Reviewer	Claudia Sagastizábal
Examiner	Jean-Charles Gilbert
Examiner	Mathurin Massias
Advisor	Franck Iutzeler
Advisor	Jérôme Malick
Guest	Claude Lemaréchal

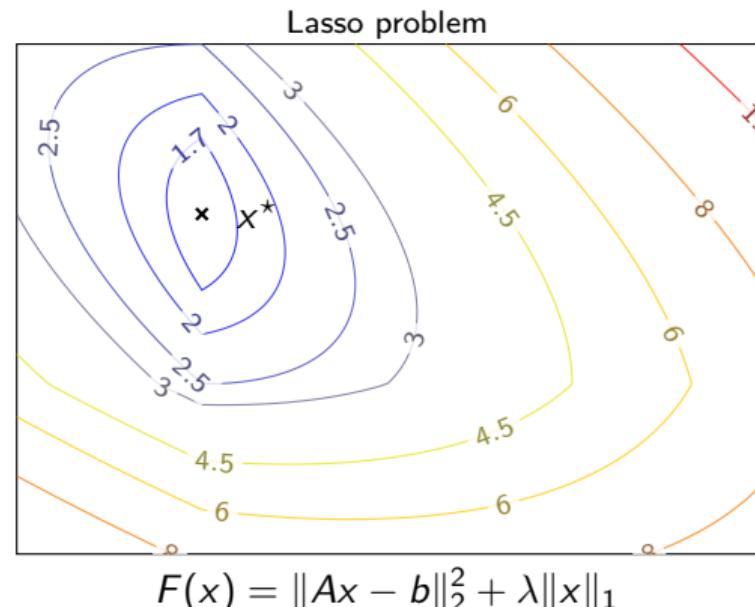
What this PhD is about

Optimization of

functions: $\min_x F(x)$



Smooth

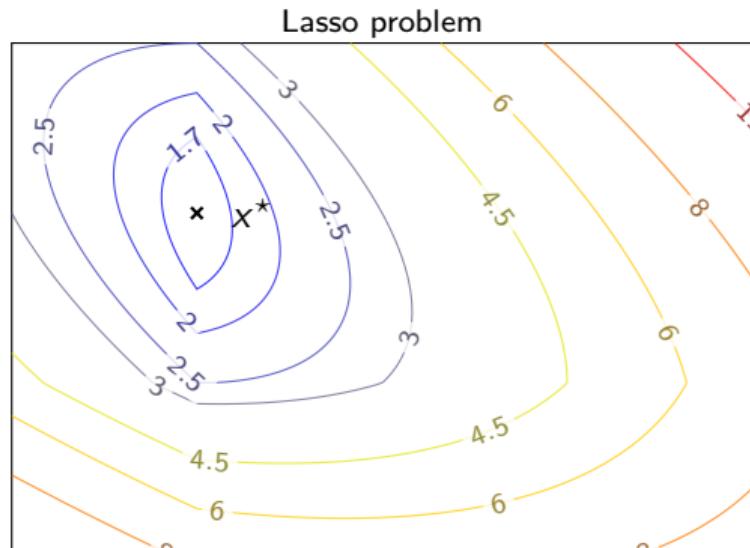


Nonsmooth

What this PhD is about

Optimization of

nonsmooth functions: $\min_x F(x)$

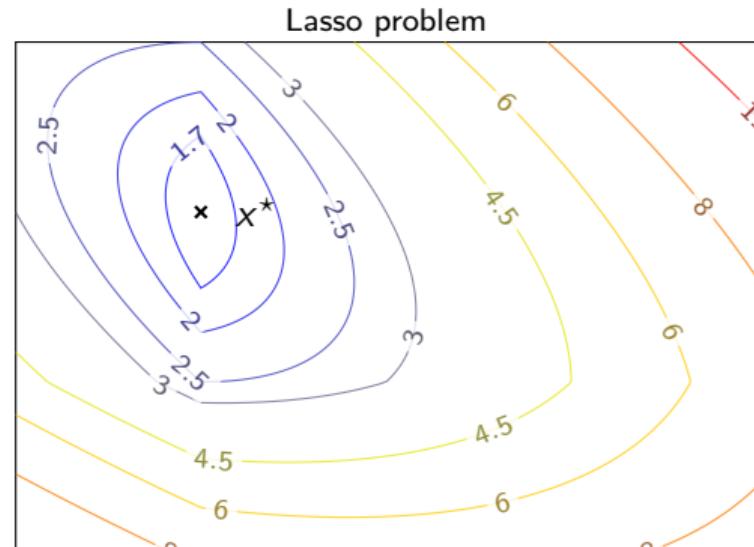


$$F(x) = \|Ax - b\|_2^2 + \lambda\|x\|_1$$

What this PhD is about

Optimization of **structured** nonsmooth functions: $\min_x F(x)$

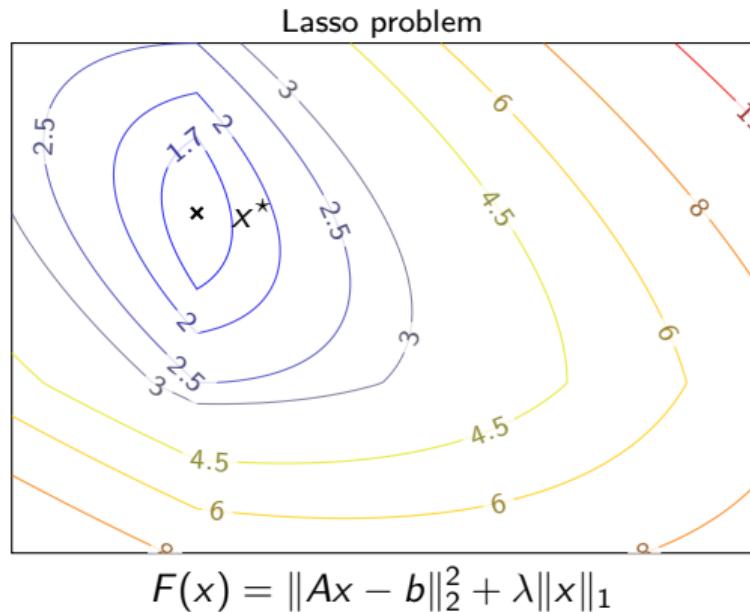
- ▷ nonsmooth points are “well-behaved”
on x & y axes



What this PhD is about

Optimization of **structured** nonsmooth functions: $\min_x F(x)$

- ▷ nonsmooth points are “well-behaved”
on x & y axes



Ok for lasso.

What about **structure in general nonsmooth functions**?

Sources of nonsmoothness 1

▷ **Implicit** nonsmoothness

$$F(x) = \sup_{u \in U} h(x, u)$$

Examples

- ▶ Robust optimization ◊ Ben-Tal *et al* '09
- ▶ Decomposition methods *e.g.*, Lagrangian relaxation, Benders decomposition ◊ BGLS '06
- ▶ Risk-averse optimization... advertising my team:
 - ▶ superquantile optim. ◊ Laguel Malick Harchaoui '22
 - ▶ distributionally robust optim. ◊ Azizian Iutzeler Malick '22

Sources of nonsmoothness 1

▷ **Implicit** nonsmoothness

$$F(x) = \sup_{u \in U} h(x, u)$$

Examples

- ▶ Robust optimization ◊ Ben-Tal *et al* '09
- ▶ Decomposition methods *e.g.*, Lagrangian relaxation, Benders decomposition ◊ BGLS '06
- ▶ Risk-averse optimization... advertising my team:
 - ▶ superquantile optim. ◊ Laguel Malick Harchaoui '22
 - ▶ distributionally robust optim. ◊ Azizian Iutzeler Malick '22
- ▶ Full **first-order** description of F near x *e.g.*, $\partial F(x)$?

Sources of nonsmoothness 1

▷ **Implicit** nonsmoothness

$$F(x) = \sup_{u \in U} h(x, u)$$

Examples

- ▶ Robust optimization ◊ Ben-Tal *et al* '09
- ▶ Decomposition methods *e.g.*, Lagrangian relaxation, Benders decomposition ◊ BGLS '06
- ▶ Risk-averse optimization... advertising my team:
 - ▶ superquantile optim. ◊ Laguel Malick Harchaoui '22
 - ▶ distributionally robust optim. ◊ Azizian Iutzeler Malick '22

▷ **Full first-order** description of F near x *e.g.*, $\partial F(x)$?

This would require the *full* set $\arg \max_{u \in U} h(x, u)$.

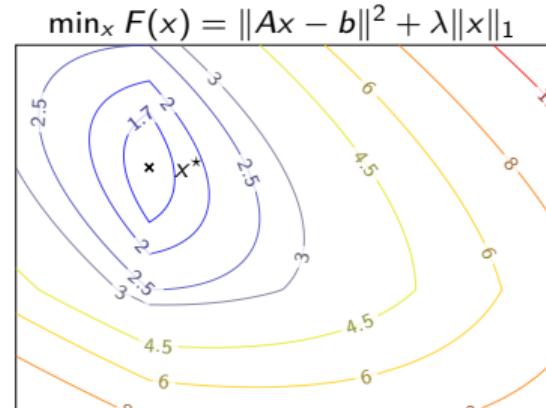
Sources of nonsmoothness 2

- ▷ **Chosen** nonsmoothness to attract minimizers

$$F(x) = f(x) + g(x), \quad \text{with } f \text{ smooth, } g \text{ nonsmooth}$$

Examples inverse problems, sparse regression e.g. lasso

- ◊ Scherzer et al '09, Vaiter et al '15



Sources of nonsmoothness 2

- ▷ **Chosen** nonsmoothness to attract minimizers

$$F(x) = f(x) + g(x), \quad \text{with } f \text{ smooth, } g \text{ nonsmooth}$$

Examples inverse problems, sparse regression e.g. lasso

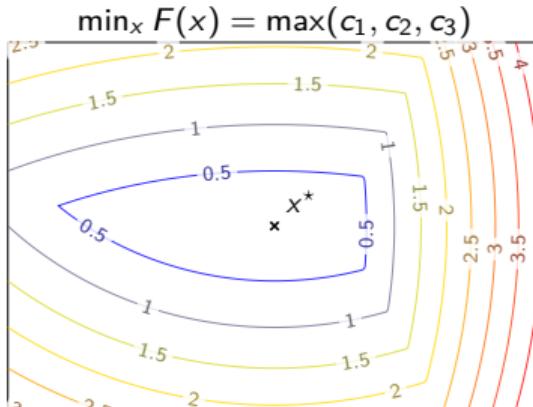
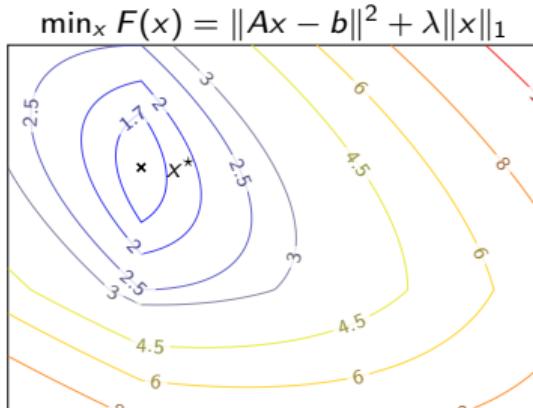
◊ Scherzer et al '09, Vaiter et al '15

- ▷ **In-between** nonsmoothness

$$F(x) = g \circ c(x), \quad \text{with } c \text{ smooth map, } g \text{ nonsmooth}$$

Examples Robust regression, optimal control, SDP

◊ Shapiro '03, Noll '05, Lewis Wright '16



Sources of nonsmoothness 2

- ▷ **Chosen** nonsmoothness to attract minimizers

$$F(x) = f(x) + g(x), \quad \text{with } f \text{ smooth, } g \text{ nonsmooth}$$

Examples inverse problems, sparse regression e.g. lasso

- ◊ Scherzer et al '09, Vaiter et al '15

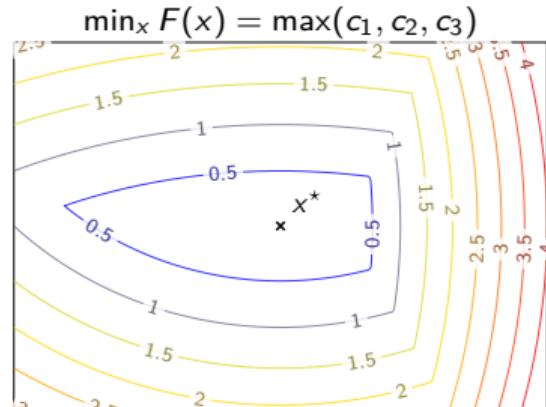
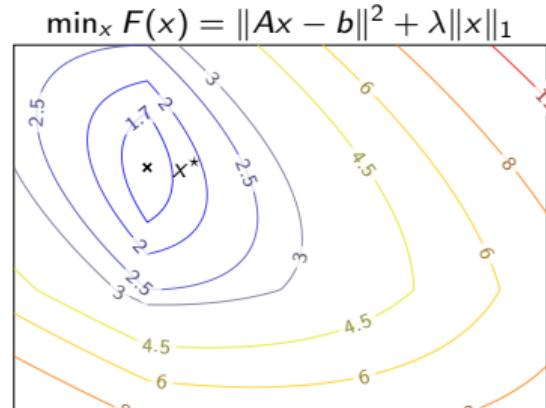
- ▷ **In-between** nonsmoothness

$$F(x) = g \circ c(x), \quad \text{with } c \text{ smooth map, } g \text{ nonsmooth}$$

Examples Robust regression, optimal control, SDP

- ◊ Shapiro '03, Noll '05, Lewis Wright '16

- ▷ In these two cases, we know



Sources of nonsmoothness 2

- ▷ **Chosen** nonsmoothness to attract minimizers

$$F(x) = f(x) + g(x), \quad \text{with } f \text{ smooth, } g \text{ nonsmooth}$$

Examples inverse problems, sparse regression e.g. lasso

◊ Scherzer et al '09, Vaiter et al '15

- ▷ **In-between** nonsmoothness

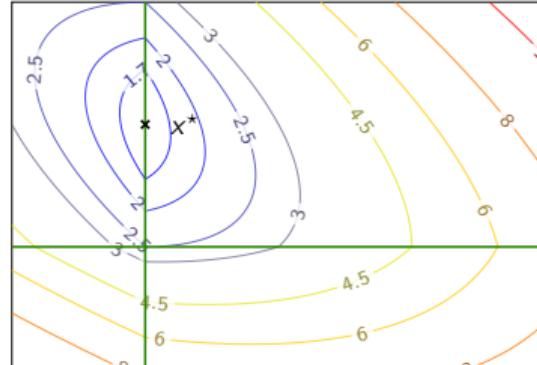
$$F(x) = g \circ c(x), \quad \text{with } c \text{ smooth map, } g \text{ nonsmooth}$$

Examples Robust regression, optimal control, SDP

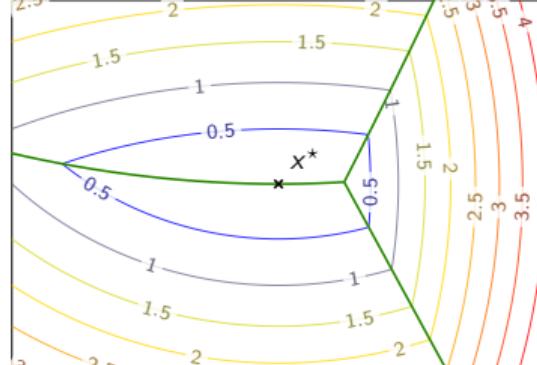
◊ Shapiro '03, Noll '05, Lewis Wright '16

- ▷ In these two cases, we know
 - ▶ where F is nonsmooth

$$\min_x F(x) = \|Ax - b\|^2 + \lambda \|x\|_1$$



$$\min_x F(x) = \max(c_1, c_2, c_3)$$



Sources of nonsmoothness 2

- ▷ **Chosen** nonsmoothness to attract minimizers

$$F(x) = f(x) + g(x), \quad \text{with } f \text{ smooth, } g \text{ nonsmooth}$$

Examples inverse problems, sparse regression e.g. lasso

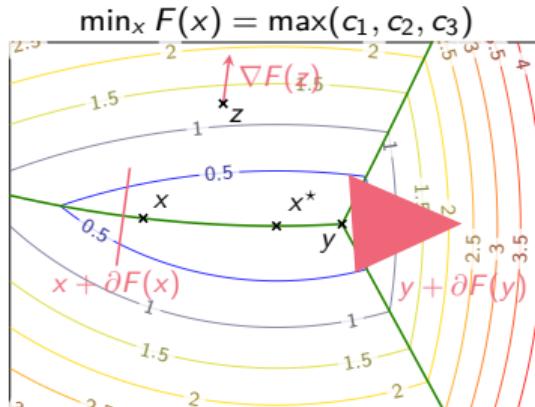
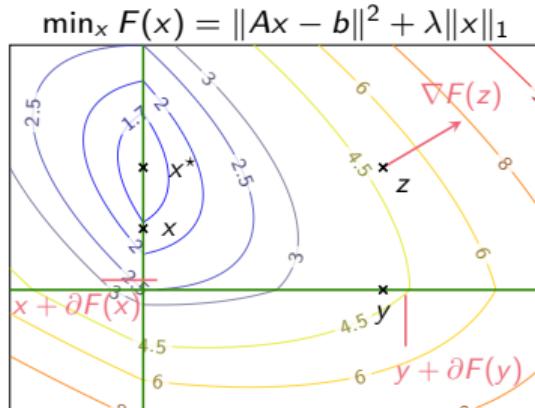
◇ Scherzer et al '09, Vaiter et al '15

- ▷ **In-between** nonsmoothness
 $F(x) = g \circ c(x)$, with c smooth map, g nonsmooth

Examples Robust regression, optimal control, SDP

◊ Shapiro '03, Noll '05, Lewis Wright '16

- ▷ In these two cases, we know
 - ▶ where F is nonsmooth
 - ▶ the full first-order description $\partial F(x)$ and more



Sources of nonsmoothness 2

- ▷ **Chosen** nonsmoothness to attract minimizers

$$F(x) = f(x) + g(x), \quad \text{with } f \text{ smooth, } g \text{ nonsmooth}$$

Examples inverse problems, sparse regression e.g. lasso

◊ Scherzer et al '09, Vaiter et al '15

- ▷ **In-between** nonsmoothness

$$F(x) = g \circ c(x), \quad \text{with } c \text{ smooth map, } g \text{ nonsmooth}$$

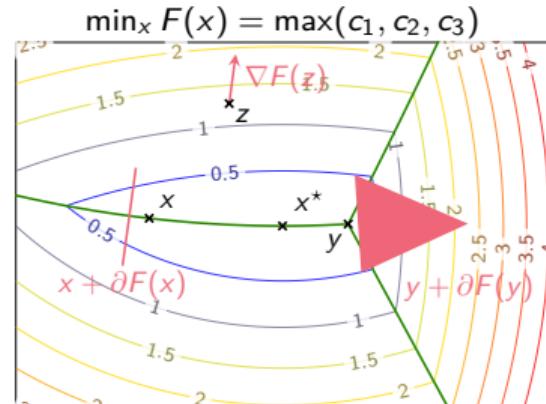
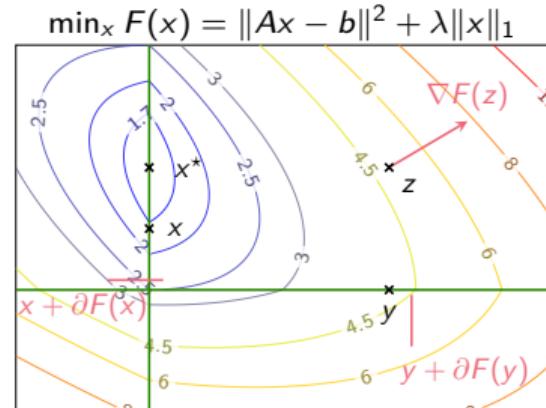
Examples Robust regression, optimal control, SDP

◊ Shapiro '03, Noll '05, Lewis Wright '16

- ▷ In these two cases, we know

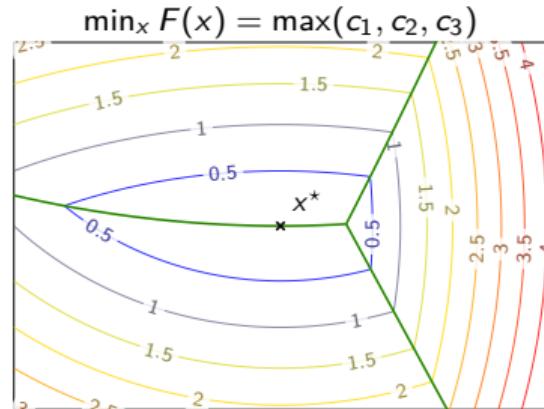
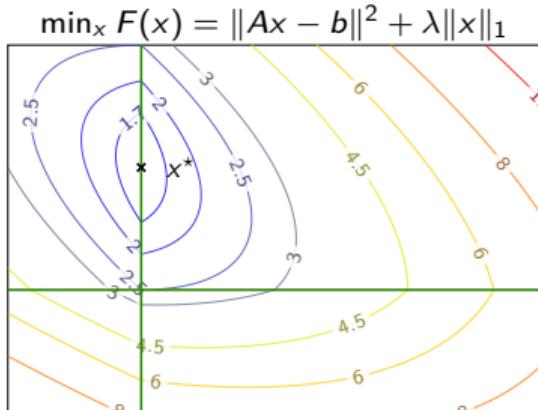
- ▶ where F is nonsmooth
- ▶ the full first-order description $\partial F(x)$ and more

Here **nonsmoothness is explicit**. Focus of this talk



Structure in nonsmoothness

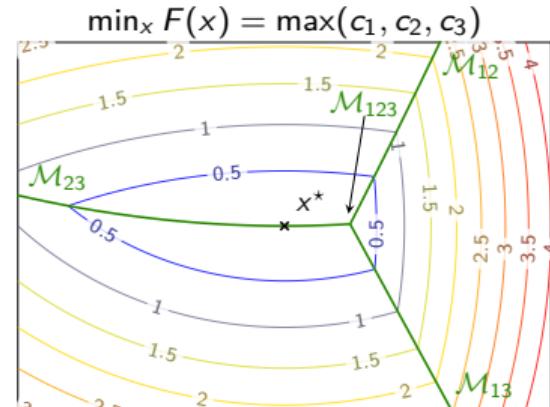
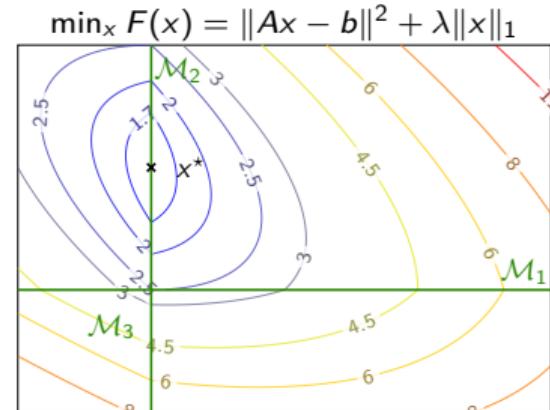
In most target applications, we **observe** that:



Structure in nonsmoothness

In most target applications, we **observe** that:

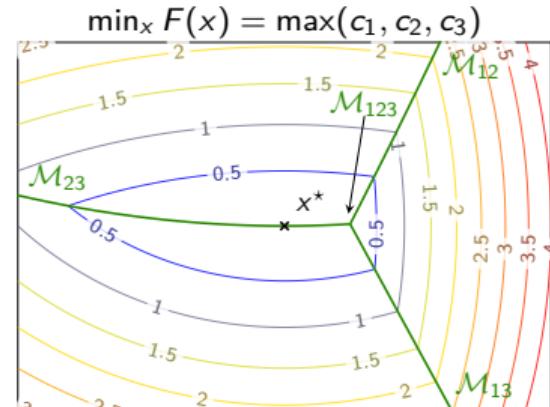
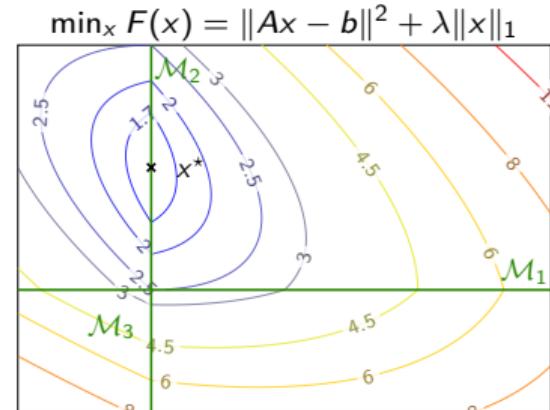
- ▶ nondiff. points organize in smooth manifolds \mathcal{M}



Structure in nonsmoothness

In most target applications, we **observe** that:

- ▶ nondiff. points organize in smooth manifolds \mathcal{M}
- ▶ locally, F is smooth along and nonsmooth across \mathcal{M}



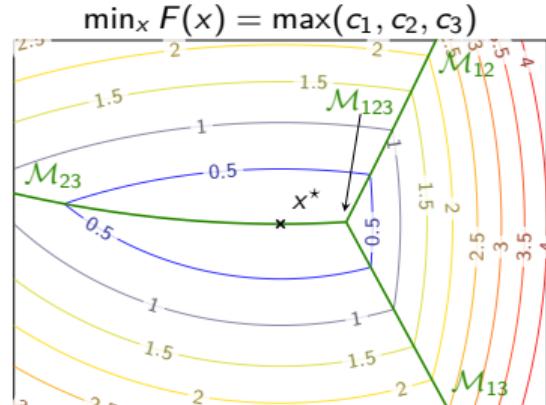
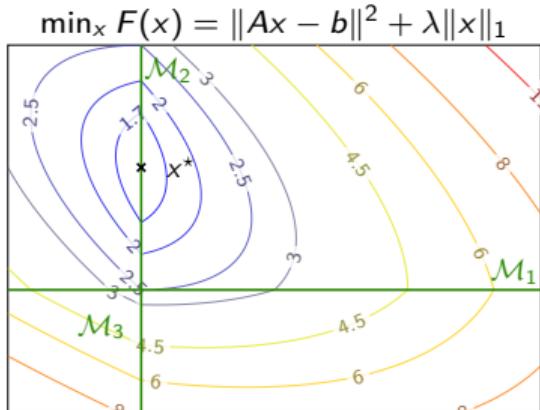
Structure in nonsmoothness

In most target applications, we **observe** that:

- ▶ nondiff. points organize in smooth manifolds \mathcal{M}
- ▶ locally, F is smooth along and nonsmooth across \mathcal{M}

These are **structure manifolds**.

◊ Lewis '02, Fadili Malick Peyré '18, Davis Drusviatsky '19



Structure in nonsmoothness

In most target applications, we **observe** that:

- ▶ nondiff. points organize in smooth manifolds \mathcal{M}
- ▶ locally, F is smooth along and nonsmooth across \mathcal{M}

These are **structure manifolds**.

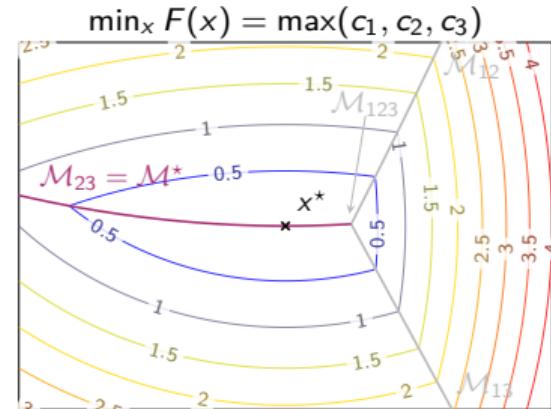
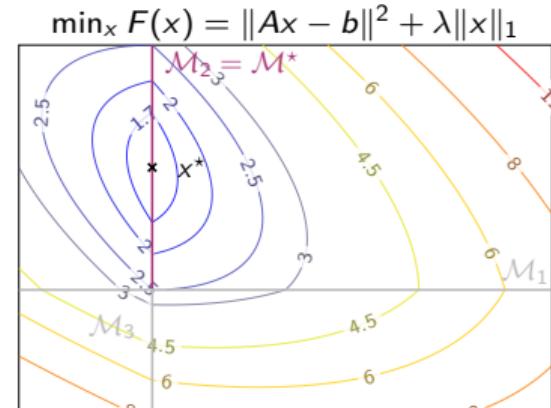
◊ Lewis '02, Fadili Malick Peyré '18, Davis Drusviatsky '19

If x^* is nonsmooth for F , there is an

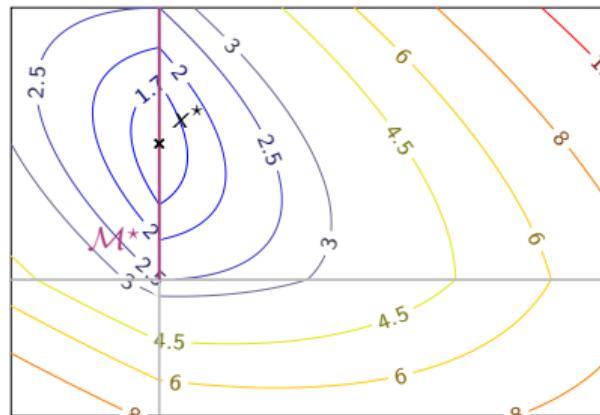
optimal manifold $\mathcal{M}^* \ni x^*$

→ This PhD explores the following question

Can structure help optimization?



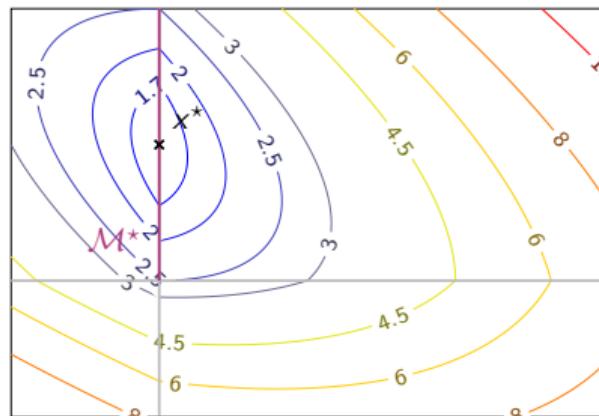
1. Detecting structure



1. Detecting structure

▷ **Implicit** nonsmoothness: $x \mapsto F(x)$, $v \in \partial F(x)$

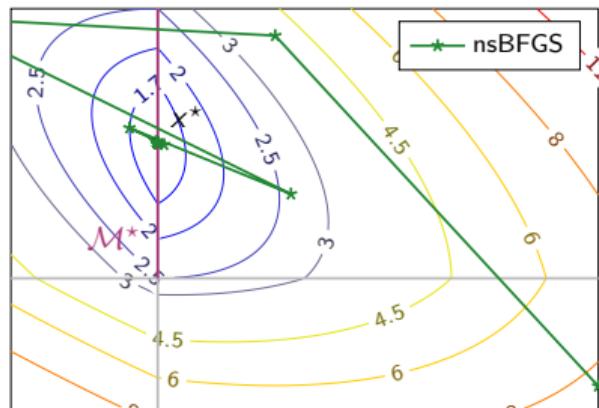
- ▶ *nonsmooth BFGS* ◇ Lewis Overton '13
- ▶ gradient sampling ◇ Lewis et al '05
- ▶ bundle methods ◇ HULL '93



1. Detecting structure

▷ **Implicit** nonsmoothness: $x \mapsto F(x)$, $v \in \partial F(x)$

- ▶ *nonsmooth BFGS* ◇ Lewis Overton '13
- ▶ gradient sampling ◇ Lewis et al '05
- ▶ bundle methods ◇ HULL '93



▷ nsBFGS: $x_k \notin \mathcal{M}^*$, oscillations
 → **no detection** of \mathcal{M}^* , suffered nonsmoothness

1. Detecting structure

▷ **Implicit** nonsmoothness: $x \mapsto F(x)$, $v \in \partial F(x)$

- ▶ *nonsmooth BFGS* ◇ Lewis Overton '13
- ▶ gradient sampling ◇ Lewis et al '05
- ▶ bundle methods ◇ HULL '93

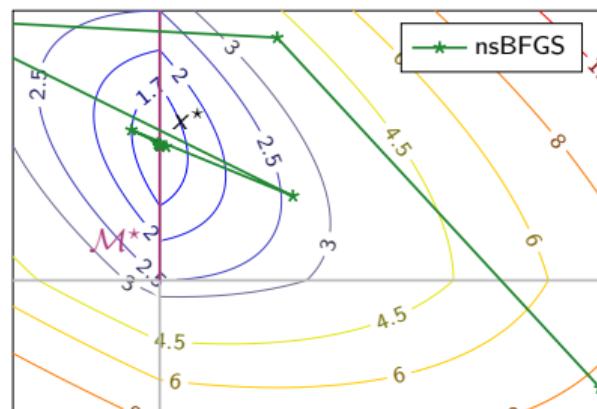
▷ **Explicit** nonsmoothness: $\text{prox}_{\gamma g}$

- ▶ *splitting methods* ◇ Combettes Pesquet, '11
- ▶ prox-linear methods ◇ Lewis Wright, '16

$$\text{prox}_{\gamma g}(y) = \arg \min_u \left\{ g(u) + \frac{1}{2\gamma} \|u - y\|^2 \right\}$$

Implicit first-order step

Closed-form for simple nonsmooth functions



▷ nsBFGS: $x_k \notin \mathcal{M}^*$, oscillations

→ **no detection** of \mathcal{M}^* , suffered nonsmoothness

1. Detecting structure

▷ **Implicit** nonsmoothness: $x \mapsto F(x)$, $v \in \partial F(x)$

- ▶ *nonsmooth BFGS* ◇ Lewis Overton '13
- ▶ gradient sampling ◇ Lewis et al '05
- ▶ bundle methods ◇ HULL '93

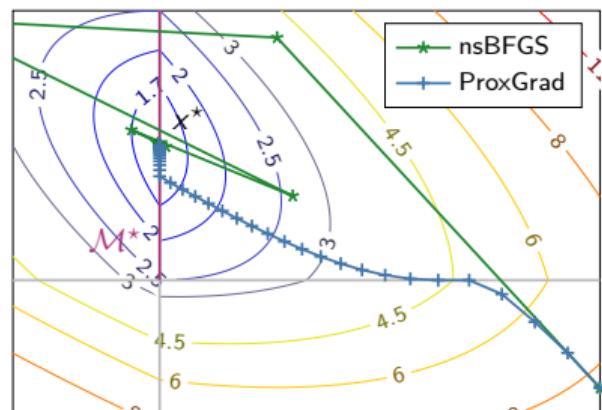
▷ **Explicit** nonsmoothness: $\text{prox}_{\gamma g}$

- ▶ *splitting methods* ◇ Combettes Pesquet, '11
- ▶ prox-linear methods ◇ Lewis Wright, '16

$$\text{prox}_{\gamma g}(y) = \arg \min_u \left\{ g(u) + \frac{1}{2\gamma} \|u - y\|^2 \right\}$$

Implicit first-order step

Closed-form for simple nonsmooth functions



▷ nsBFGS: $x_k \notin \mathcal{M}^*$, oscillations

→ **no detection** of \mathcal{M}^* , suffered nonsmoothness

▷ ProxGrad: after some time $x_k \in \mathcal{M}^*$

→ **detection** of \mathcal{M}^* , but **no exploitation**

2. Exploiting structure

- With knowledge of \mathcal{M}^* , nonsmooth min. turns into smooth constrained opt.

$$\min_{x \in \mathbb{R}^n} F(x) \quad \xrightarrow{\text{with } \mathcal{M}^*} \quad \min_{x \in \mathcal{M}^*} F|_{\mathcal{M}^*}(x)$$

2. Exploiting structure

- With knowledge of \mathcal{M}^* , nonsmooth min. turns into smooth constrained opt.

$$\min_{x \in \mathbb{R}^n} F(x) \quad \xrightarrow{\text{with } \mathcal{M}^*} \quad \min_{x \in \mathcal{M}^*} F|_{\mathcal{M}^*}(x)$$

- Tune splitting method parameters → linear rate ◇ Liang Fadili Peyré '17
- \mathcal{VU} -algorithm → superlinear steps ◇ Mifflin Sagastizábal '05
 - Approximate Newton on smooth directions, approximate prox

These algorithms **exploit** knowledge of \mathcal{M}^*

2. Exploiting structure

- With knowledge of \mathcal{M}^* , nonsmooth min. turns into smooth constrained opt.

$$\min_{x \in \mathbb{R}^n} F(x) \quad \xrightarrow{\text{with } \mathcal{M}^*} \quad \min_{x \in \mathcal{M}^*} F|_{\mathcal{M}^*}(x)$$

- Tune splitting method parameters → linear rate ◇ Liang Fadili Peyré '17
- \mathcal{VU} -algorithm → superlinear steps ◇ Mifflin Sagastizábal '05
Approximate Newton on smooth directions, approximate prox

These algorithms **exploit** knowledge of \mathcal{M}^*

Problem: \mathcal{M}^* is not known, only approximated numerically
Except in specific cases ◇ Daniilidis Sagastizábal Solodov '09

This thesis: detecting & exploiting structure

This thesis: detecting & exploiting structure

- ▷ Smooth reduced problem

$$\min_{x \in \mathbb{R}^n} F(x) \quad \xrightarrow{\text{with } \mathcal{M}^*} \quad \min_{x \in \mathcal{M}^*} F|_{\mathcal{M}^*}(x)$$

If \mathcal{M}^* is known, we can use **fast Newton-type methods** relative to \mathcal{M}^*

This thesis: detecting & exploiting structure

- ▷ Smooth reduced problem

$$\min_{x \in \mathbb{R}^n} F(x) \quad \xrightarrow{\text{with } \mathcal{M}^*} \quad \min_{x \in \mathcal{M}^*} F|_{\mathcal{M}^*}(x)$$

If \mathcal{M}^* is known, we can use **fast Newton-type methods** relative to \mathcal{M}^*

Challenge: \mathcal{M}^* is never known; but near point x_k , there is relevant structure \mathcal{M}_k

This thesis: detecting & exploiting structure

- ▷ Smooth reduced problem

$$\min_{x \in \mathbb{R}^n} F(x) \quad \xrightarrow{\text{with } \mathcal{M}^*} \quad \min_{x \in \mathcal{M}^*} F|_{\mathcal{M}^*}(x)$$

If \mathcal{M}^* is known, we can use **fast Newton-type methods** relative to \mathcal{M}^*

Challenge: \mathcal{M}^* is never known; but near point x_k , there is relevant structure \mathcal{M}_k

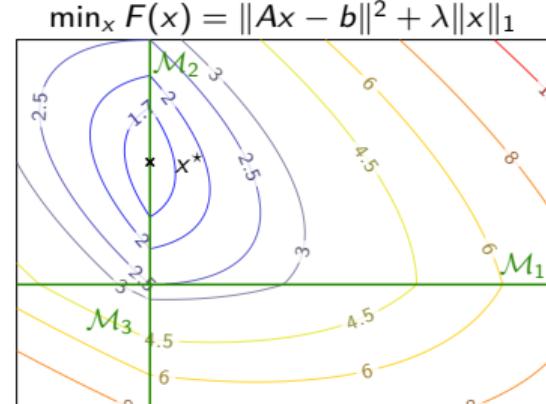
- ▷ We must use structure **adaptively**, which raises two difficulties:

- ▶ How to *provably detect* relevant structure near a point x_k ? $\rightarrow \text{prox}_{\gamma g}$
- ▶ How to **exploit** the detected structure \mathcal{M}_k ? \rightarrow Newton method

Outline

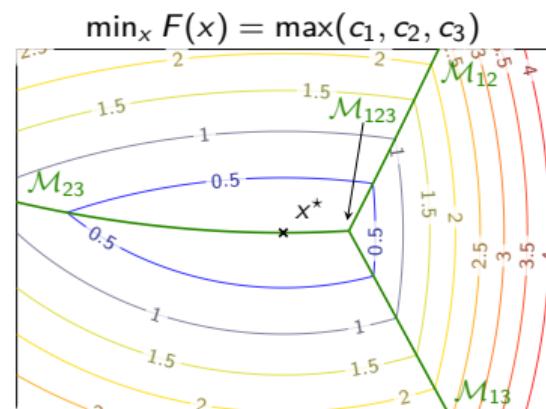
1. Additive nonsmoothness: $\min_x f(x) + g(x)$

Iterates naturally live on \mathcal{M} , feasible optim. methods



2. Composite nonsmoothness: $\min_x g \circ c(x)$

Iterates don't live on \mathcal{M} , infeasible optim. methods



Introduction
○○○○○○○

Additive nonsmoothness $f + g$
●○○○○○○○○○○

Composite nonsmoothness $g \circ c$
○○○○○○○○○○○○○○○○

Conclusion
○○○○○

Outline

Introduction

Additive nonsmoothness $f + g$

Composite nonsmoothness $g \circ c$

Conclusion

Additive problems

$$\text{Find } x^* \in \arg \min_{x \in \mathbb{R}^n} F(x) = \begin{array}{c} f(x) \\ \text{smooth} \end{array} + \begin{array}{c} g(x) \\ \text{non smooth} \end{array}$$

Example:

In inverse problems / learning, g is **chosen** to encode **prior knowledge**:

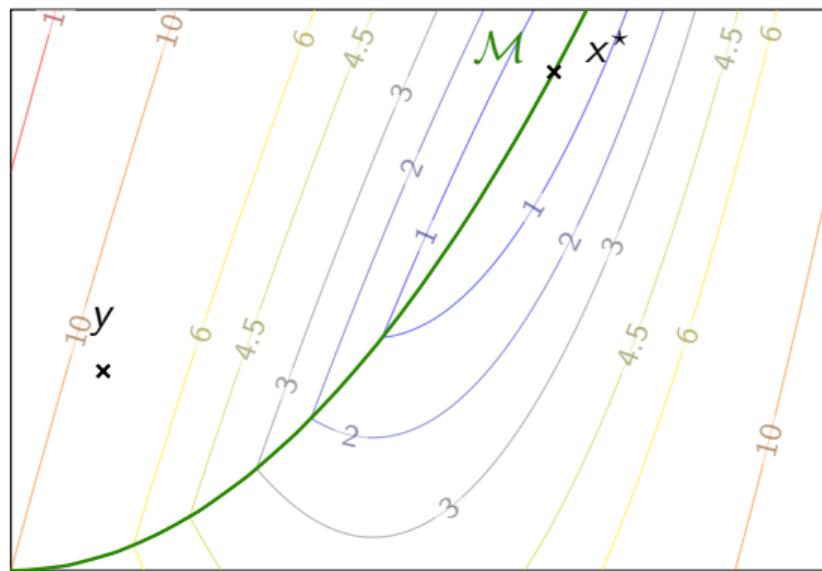
- ▶ x^* sparse $\rightarrow g(x) \propto \|x\|_1 = \sum_{i=1}^n |x_i|$
- ▶ X^* low-rank $\rightarrow g(X) \propto \|X\|_* = \sum_{i=1}^{\text{rank}(X)} \sigma_i(X)$

◊ Candès Romberg Tao '06; Candès Recht '13; Vaiter Peyré Fadili '15

Back on the proximal gradient

- ▷ Proximal gradient algorithm: iterate

$$x = \text{prox}_{\gamma g}(y - \gamma \nabla f(y))$$



$$\min_{x \in \mathbb{R}^2} F(x) = 10(1 - x_1)^2 + 5|x_1^2 - x_2|$$

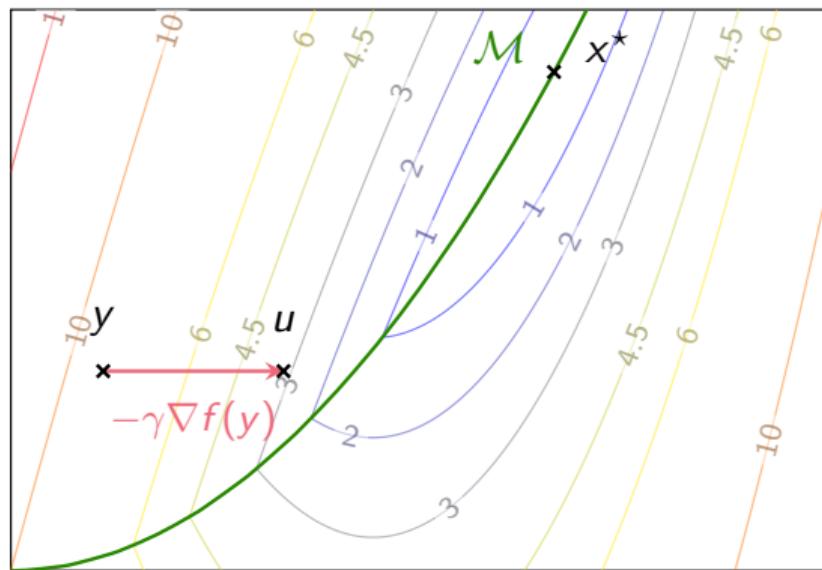
Back on the proximal gradient

- ▷ Proximal gradient algorithm: iterate

$$x = \text{prox}_{\gamma g}(y - \gamma \nabla f(y))$$

1. Explicit step on f

$$u = y - \gamma \nabla f(y)$$



$$\min_{x \in \mathbb{R}^2} F(x) = 10(1 - x_1)^2 + 5|x_1^2 - x_2|$$

Back on the proximal gradient

- Proximal gradient algorithm: iterate

$$x = \text{prox}_{\gamma g}(y - \gamma \nabla f(y))$$

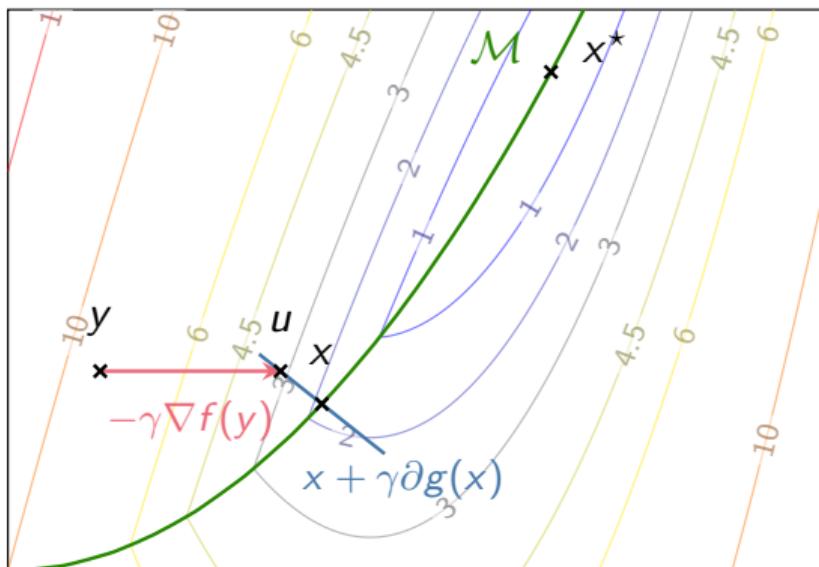
- Explicit step on f

$$u = y - \gamma \nabla f(y)$$

- Implicit step on g

$$\begin{aligned} x &= \text{prox}_{\gamma g}(u) \\ \Leftrightarrow u &\in x + \gamma \partial g(x) \end{aligned}$$

→ $\text{prox}_{\gamma g}$ can send points to \mathcal{M}



$$\min_{x \in \mathbb{R}^2} F(x) = 10(1 - x_1)^2 + 5|x_1^2 - x_2|$$

Back on the proximal gradient

- Proximal gradient algorithm: iterate

$$x = \text{prox}_{\gamma g}(y - \gamma \nabla f(y))$$

- Explicit step on f

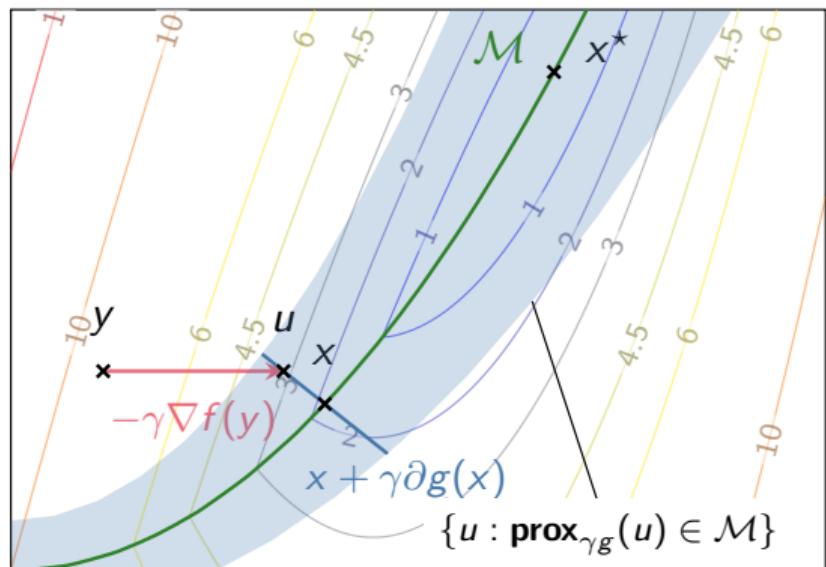
$$u = y - \gamma \nabla f(y)$$

- Implicit step on g

$$\begin{aligned} x &= \text{prox}_{\gamma g}(u) \\ \Leftrightarrow u &\in x + \gamma \partial g(x) \end{aligned}$$

→ $\text{prox}_{\gamma g}$ can send points to \mathcal{M}

More precisely...



$$\min_{x \in \mathbb{R}^2} F(x) = 10(1 - x_1)^2 + 5|x_1^2 - x_2|$$

Back on the proximal gradient

- Proximal gradient algorithm: iterate

$$x = \text{prox}_{\gamma g}(y - \gamma \nabla f(y))$$

- Explicit step on f

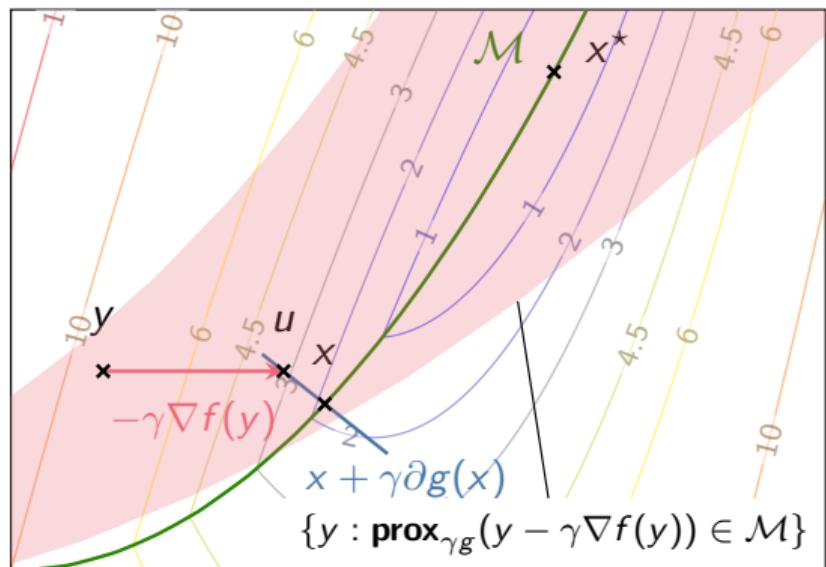
$$u = y - \gamma \nabla f(y)$$

- Implicit step on g

$$\begin{aligned} x &= \text{prox}_{\gamma g}(u) \\ \Leftrightarrow u &\in x + \gamma \partial g(x) \end{aligned}$$

→ $\text{prox}_{\gamma g}$ can send points to \mathcal{M}

More precisely...



$$\min_{x \in \mathbb{R}^2} F(x) = 10(1 - x_1)^2 + 5|x_1^2 - x_2|$$

Proximal gradient: a structure detector

- ▷ Near manifolds, the proximal gradient sends points to the manifold, smoothly:

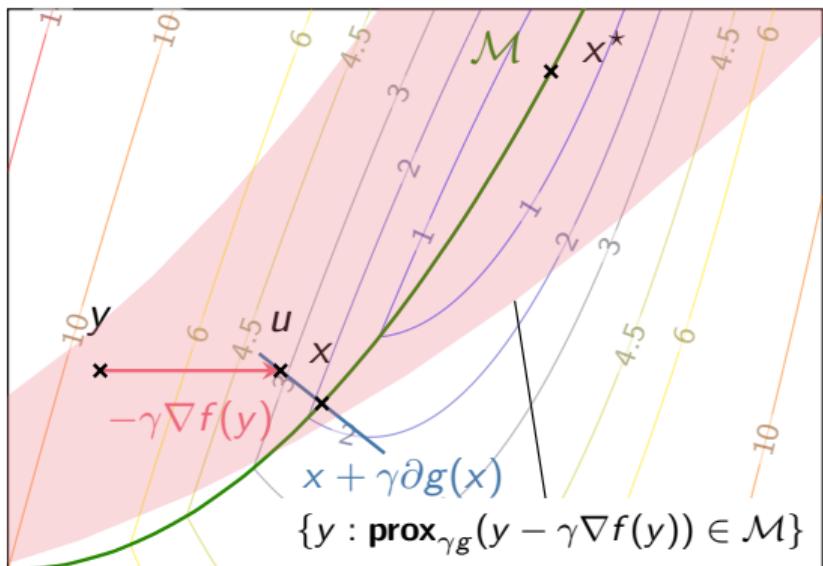
Theorem

Take y and $x = \text{prox}_{\gamma g}(y - \gamma \nabla f(y))$ s.t.

- ▶ $x \in \mathcal{M}$ partial smoothness
- ▶ $y - \gamma \nabla f(y) - x \in \text{ri } \gamma \partial g(x)$

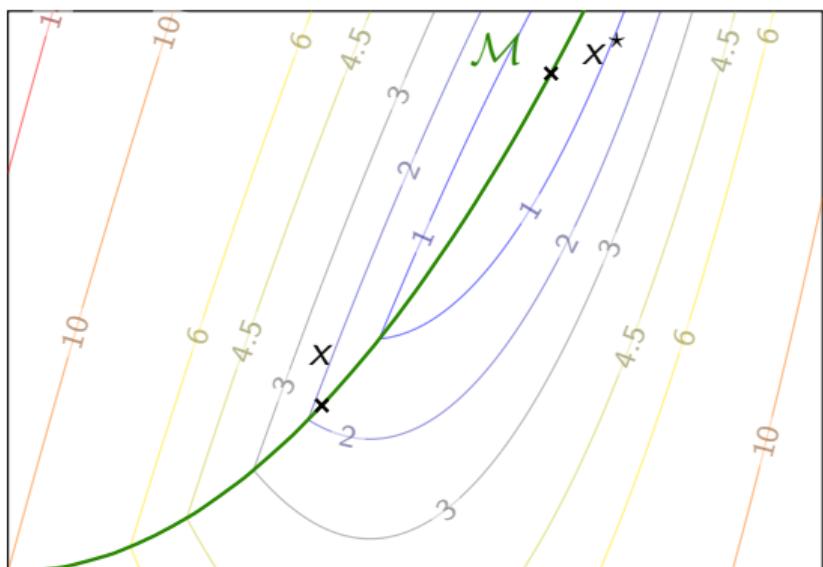
then, on a neighborhood \mathcal{N}_y of y the proximal gradient is

- ▶ \mathcal{M} -valued
- ▶ a \mathcal{C}^1 operator



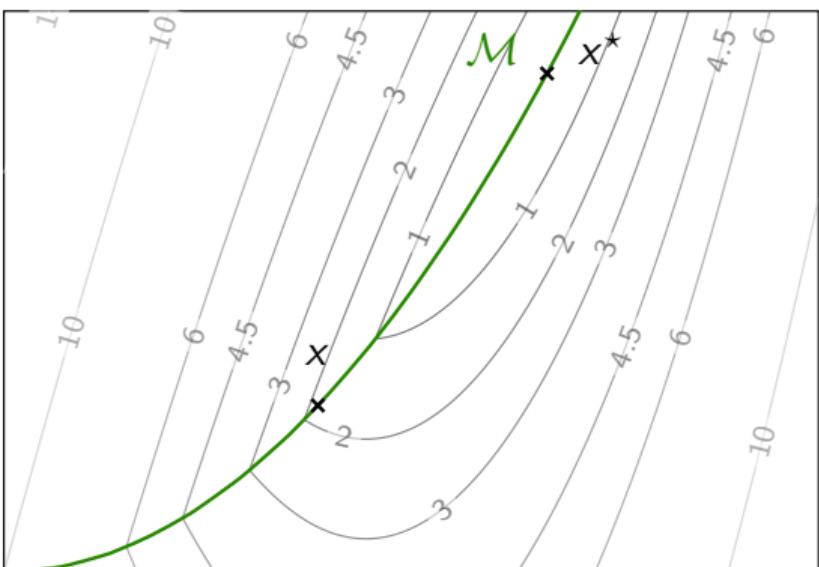
$$\min_{x \in \mathbb{R}^2} F(x) = 10(1 - x_1)^2 + 5|x_1^2 - x_2|$$

Riemannian optimization, in a nutshell



$$\min_{x \in \mathbb{R}^2} 10(1 - x_1)^2 + 5|x_1^2 - x_2| \text{ s.t. } x \in \mathcal{M}$$

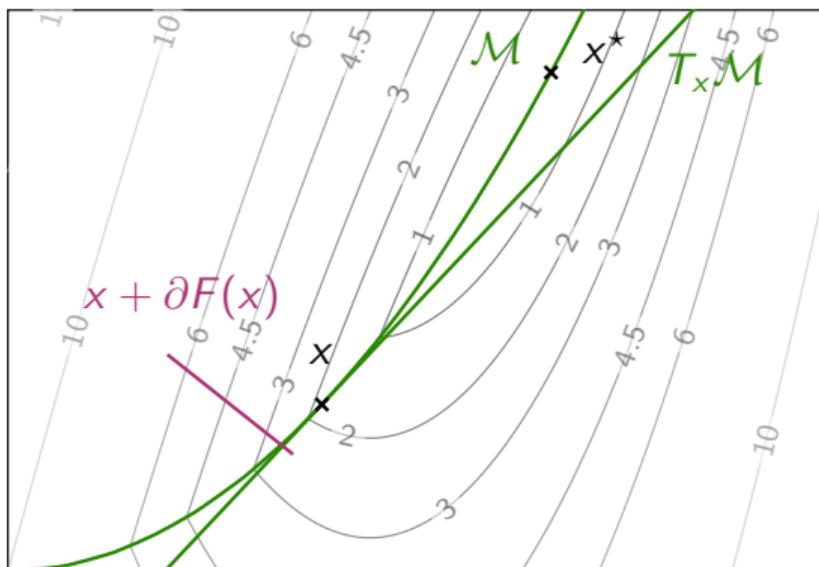
Riemannian optimization, in a nutshell



$$\min_{x \in \mathbb{R}^2} 10(1 - x_1)^2 + 5|x_1^2 - x_2| \text{ s.t. } x \in \mathcal{M}$$

Riemannian optimization, in a nutshell

Elementary tools:

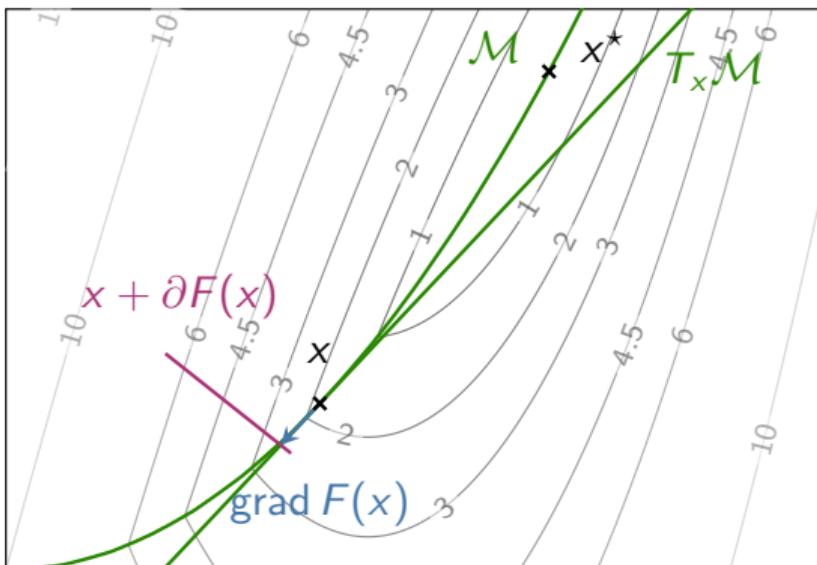


$$\min_{x \in \mathbb{R}^2} 10(1 - x_1)^2 + 5|x_1^2 - x_2| \text{ s.t. } x \in \mathcal{M}$$

Riemannian optimization, in a nutshell

Elementary tools:

- ▶ Riemannian gradient, Hessian



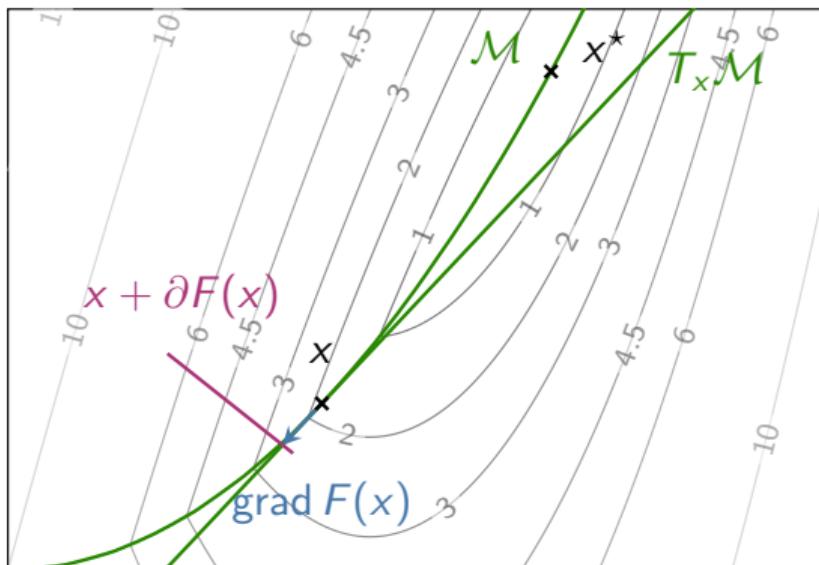
$$\min_{x \in \mathbb{R}^2} 10(1 - x_1)^2 + 5|x_1^2 - x_2| \text{ s.t. } x \in \mathcal{M}$$

Riemannian optimization, in a nutshell

Elementary tools:

- ▶ Riemannian gradient, Hessian

Typical Riemannian step



$$\min_{x \in \mathbb{R}^2} 10(1 - x_1)^2 + 5|x_1^2 - x_2| \text{ s.t. } x \in \mathcal{M}$$

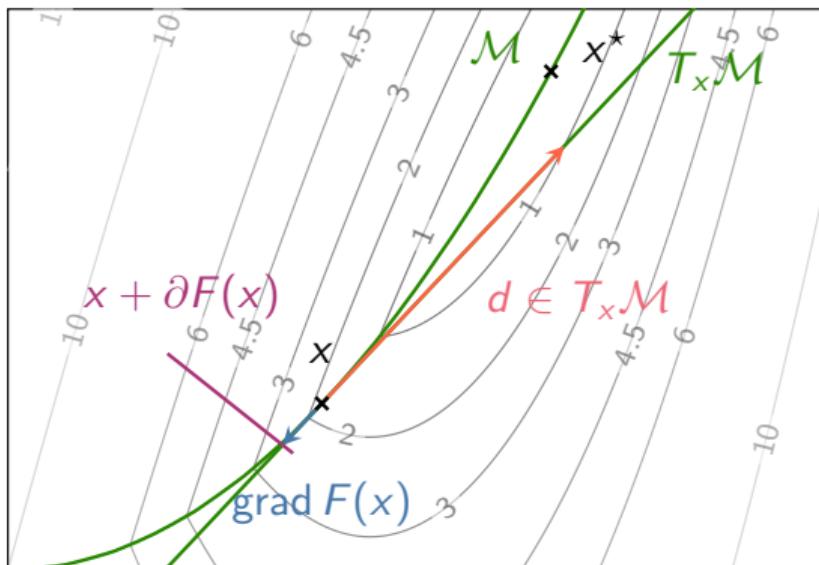
Riemannian optimization, in a nutshell

Elementary tools:

- ▶ Riemannian gradient, Hessian

Typical Riemannian step

- ▶ find $d \in T_x \mathcal{M}$ e.g., $d = -\alpha \text{grad } F(x)$



$$\min_{x \in \mathbb{R}^2} 10(1 - x_1)^2 + 5|x_1^2 - x_2| \text{ s.t. } x \in \mathcal{M}$$

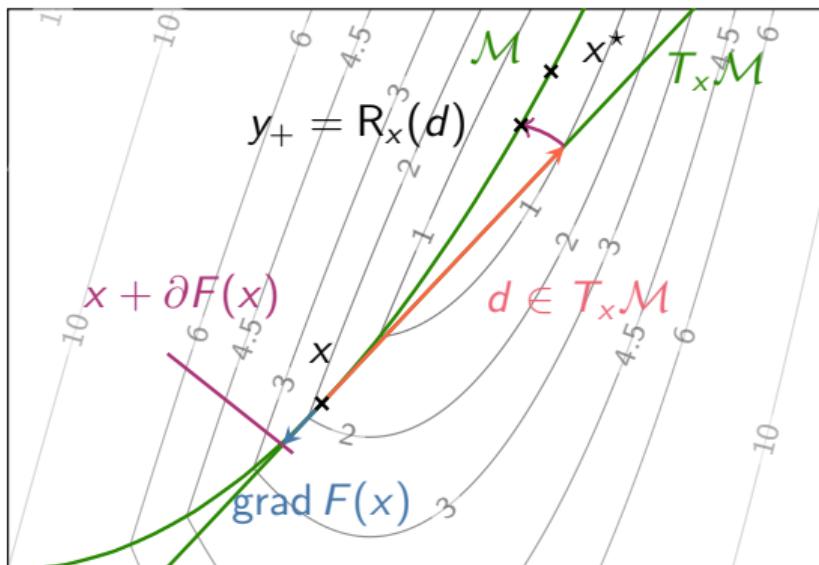
Riemannian optimization, in a nutshell

Elementary tools:

- ▶ Riemannian gradient, Hessian

Typical Riemannian step

- ▶ find $d \in T_x \mathcal{M}$ e.g., $d = -\alpha \text{grad } F(x)$
- ▶ send d to \mathcal{M} e.g., retraction $R_x(d)$



$$\min_{x \in \mathbb{R}^2} 10(1 - x_1)^2 + 5|x_1^2 - x_2| \text{ s.t. } x \in \mathcal{M}$$

Riemannian optimization, in a nutshell

Elementary tools:

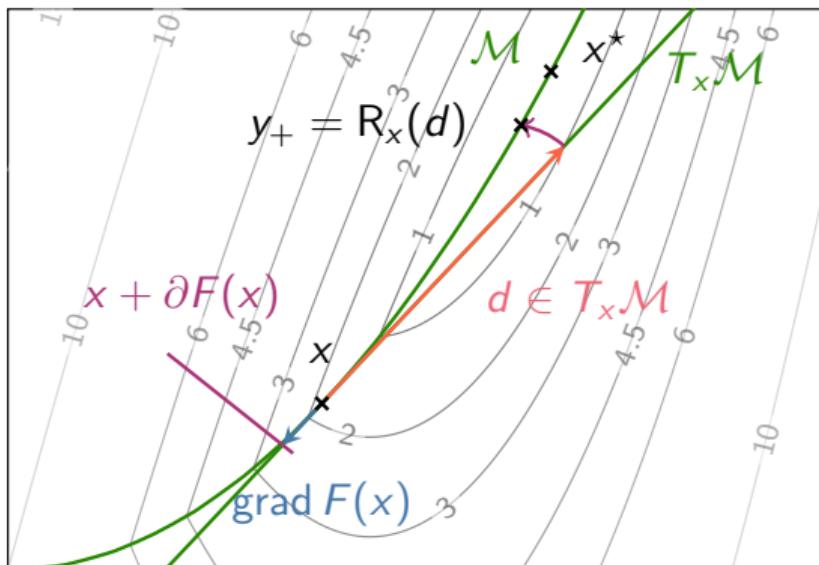
- ▶ Riemannian gradient, Hessian

Typical Riemannian step

- ▶ find $d \in T_x \mathcal{M}$ e.g., $d = -\alpha \text{grad } F(x)$
- ▶ send d to \mathcal{M} e.g., retraction $R_x(d)$

→ Many methods of smooth optim carry over manifold optimization.

◊ Absil et al '09, Boumal '22



$$\min_{x \in \mathbb{R}^2} 10(1 - x_1)^2 + 5|x_1^2 - x_2| \text{ s.t. } x \in \mathcal{M}$$

Proposed algorithm

Iteration k :

- ▷ Compute point $x_k = \text{prox}_{\gamma g}(y_{k-1} - \gamma \nabla f(y_{k-1}))$ and manifold $\mathcal{M}_k \ni x_k$
- ▷ Compute y_k from point x_k by a manifold update on \mathcal{M}_k e.g., Riemannian Newton step

Proposed algorithm

Iteration k :

- ▷ Compute point $x_k = \text{prox}_{\gamma g}(y_{k-1} - \gamma \nabla f(y_{k-1}))$ and manifold $\mathcal{M}_k \ni x_k$
- ▷ Compute y_k from point x_k by a manifold update on \mathcal{M}_k e.g., Riemannian Newton step

Theorem

- ▷ If $\gamma < 1/L$ and the manifold update decreases function value, then any limit point \bar{x} of (x_k) is a critical point: $0 \in \partial F(\bar{x})$

Proposed algorithm

Iteration k :

- ▷ Compute point $x_k = \text{prox}_{\gamma g}(y_{k-1} - \gamma \nabla f(y_{k-1}))$ and manifold $\mathcal{M}_k \ni x_k$
- ▷ Compute y_k from point x_k by a manifold update on \mathcal{M}_k e.g., Riemannian Newton step

Theorem

- ▷ If $\gamma < 1/L$ and the manifold update decreases function value, then any limit point \bar{x} of (x_k) is a critical point: $0 \in \partial F(\bar{x})$
- ▷ Assume that a Riemannian Newton method is used, and that one limit point x^* is such that
 - ▶ g has structure \mathcal{M}^* at x^*
 - ▶ $0 \in \text{ri } \partial F(x^*)$, $\text{Hess}_{\mathcal{M}^*} F(x^*) \succ 0$ and $\text{Hess}_{\mathcal{M}^*}$ is locally Lipschitz around x^*

Then, after some finite time

- ▶ $x_k \in \mathcal{M}^*$
- ▶ x_k converges to x^* at a **quadratic rate**: $\|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2$

Illustration: toy example

$$\min_{x \in \mathbb{R}^2} \underbrace{2x_1^2 + x_2^2}_{f(x)} + \underbrace{|x_1^2 - x_2|}_{g(x)}$$

- +— Proximal Gradient
- *— Accel. Proximal Gradient
- ⊕— Alt. Newton

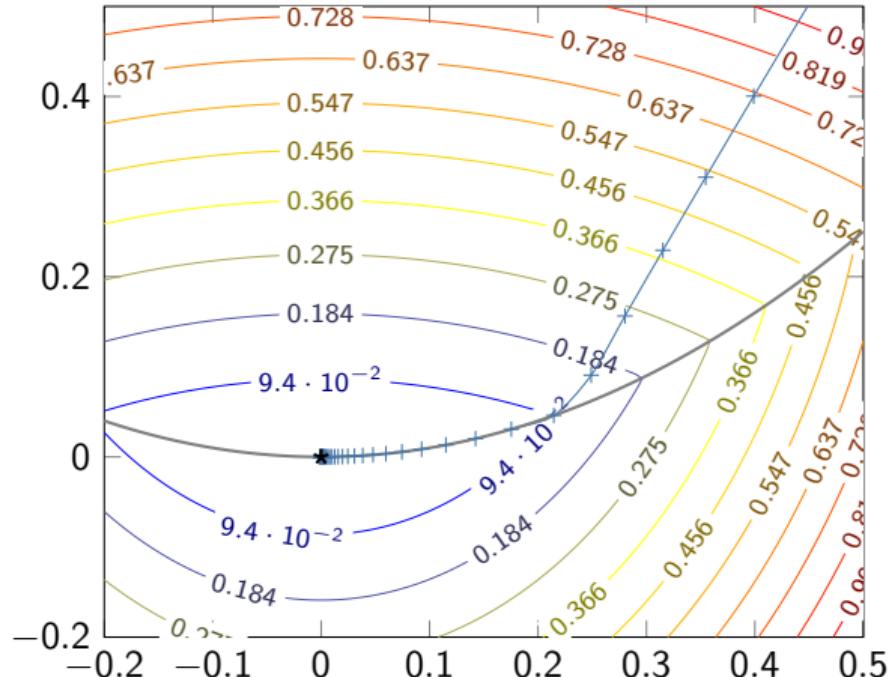


Illustration: toy example

$$\min_{x \in \mathbb{R}^2} \underbrace{2x_1^2 + x_2^2}_{f(x)} + \underbrace{|x_1^2 - x_2|}_{g(x)}$$

- +— Proximal Gradient
- *— Accel. Proximal Gradient
- ⊕— Alt. Newton

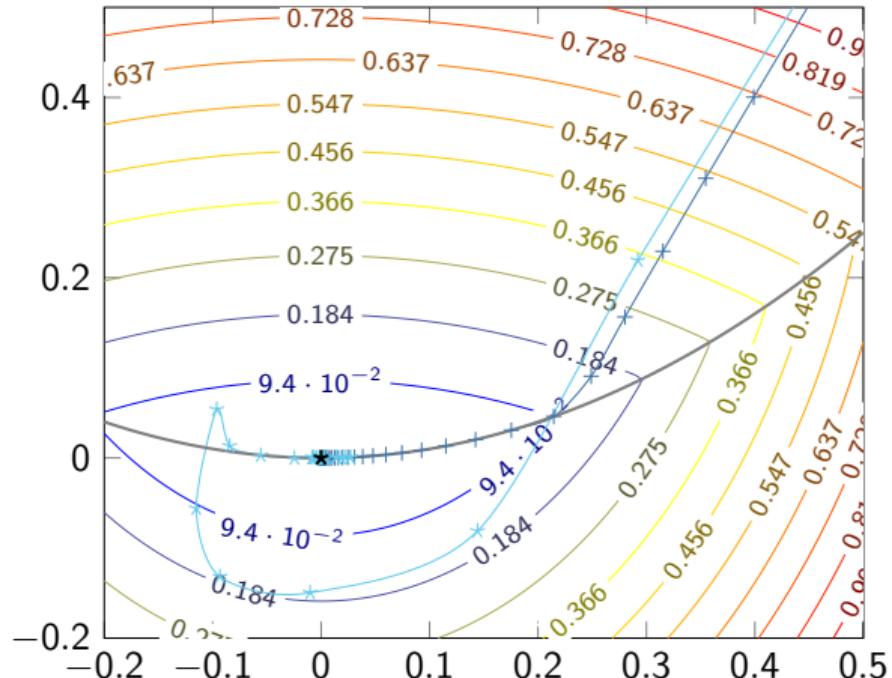


Illustration: toy example

$$\min_{x \in \mathbb{R}^2} \underbrace{2x_1^2 + x_2^2}_{f(x)} + \underbrace{|x_1^2 - x_2|}_{g(x)}$$

- +— Proximal Gradient
- *— Accel. Proximal Gradient
- ⊕— Alt. Newton

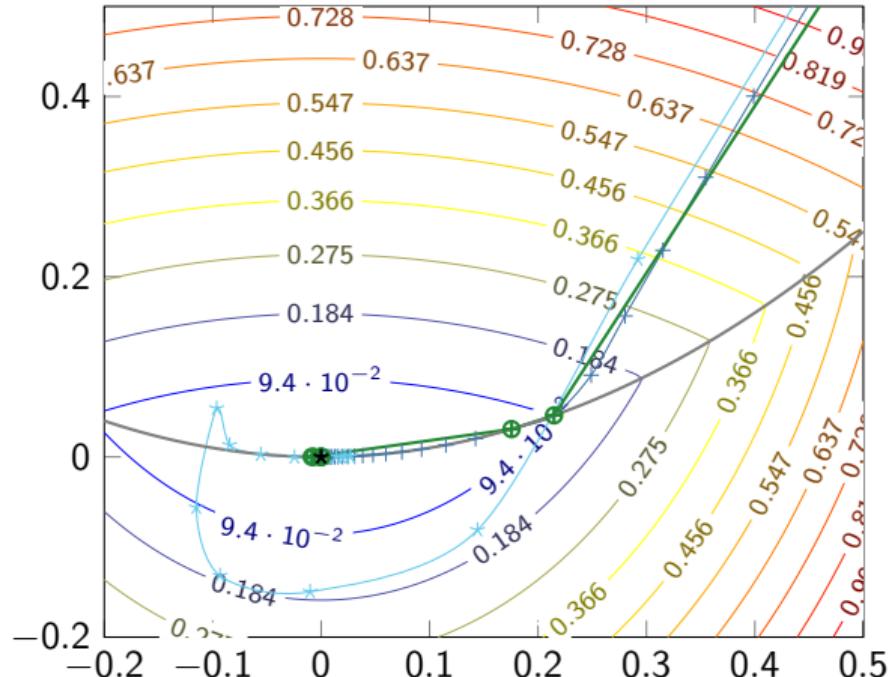


Illustration 1: logistic regression

$$\min_{x \in \mathbb{R}^n} \sum_{i=1}^m \log(y_i \sigma(\langle A_i, x \rangle)) + \lambda \|x\|_1$$

In this instance, $n = m = 4000$

$$\dim(\mathcal{M}^*) = 249$$

- +— Proximal Gradient
- *— Accel. Proximal Gradient
- ⊕— Alt. Newton
- △— Alt. Truncated Newton

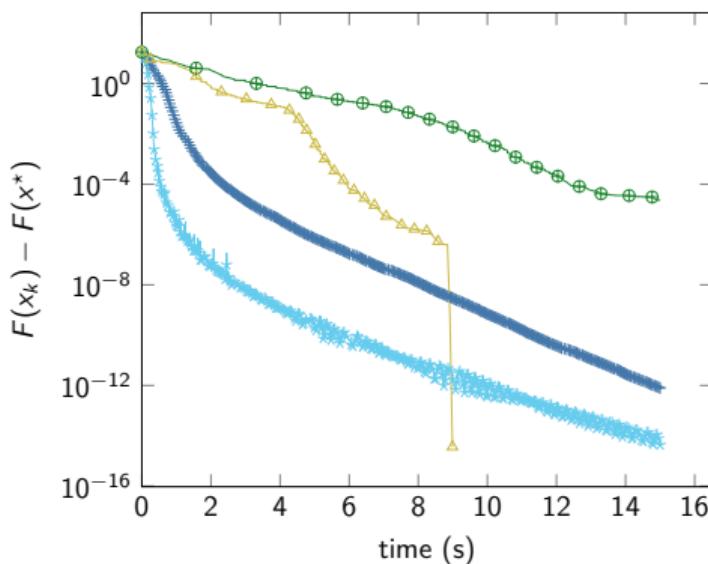


Illustration 1: logistic regression

$$\min_{x \in \mathbb{R}^n} \sum_{i=1}^m \log(y_i \sigma(\langle A_i, x \rangle)) + \lambda \|x\|_1$$

In this instance, $n = m = 4000$

$$\dim(\mathcal{M}^*) = 249$$

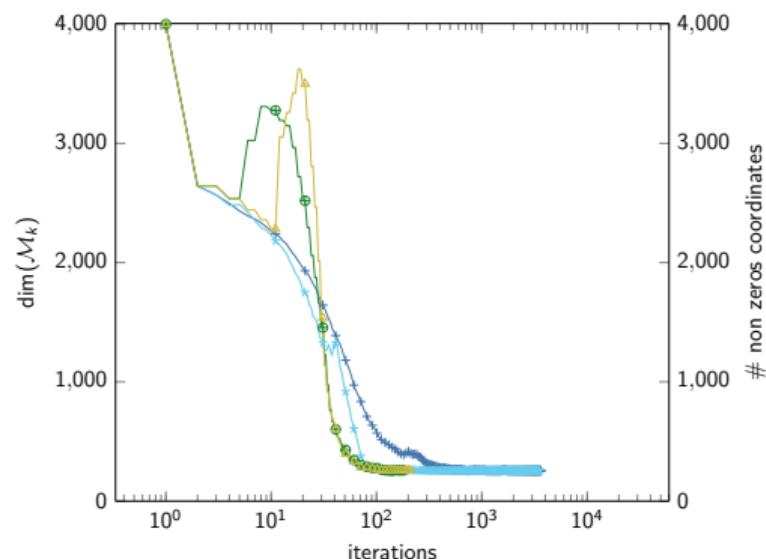
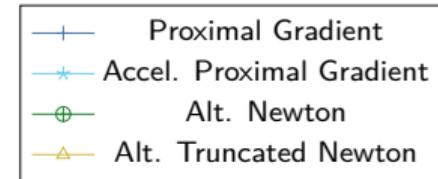
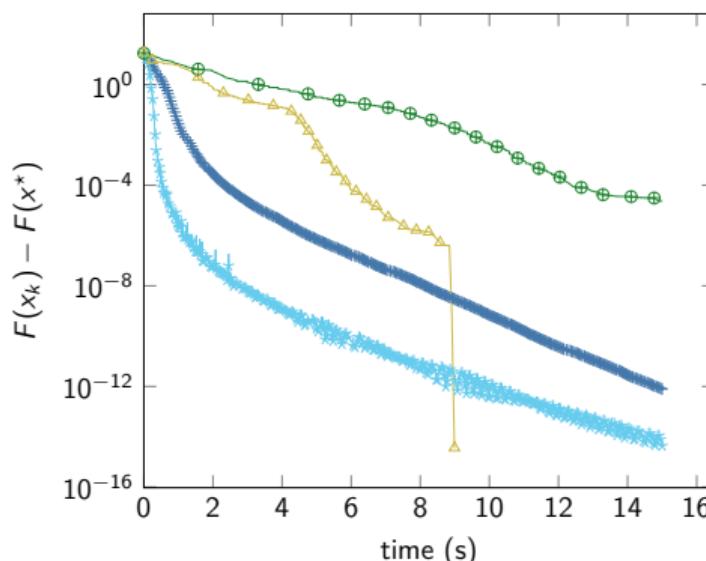


Illustration 2: tracenorm regression

$$\min_{X \in \mathbb{R}^{10 \times 12}} \sum_{i=1}^m (\langle A_i, X \rangle - y_i)^2 + \lambda \|X\|_*$$

In this instance, $m = 60$

$$\mathcal{M}^* = \{X : \text{rank}(X) = 6\}$$

- +— Proximal Gradient
- *— Accel. Proximal Gradient
- ⊕— Alt. Newton
- △— Alt. Truncated Newton

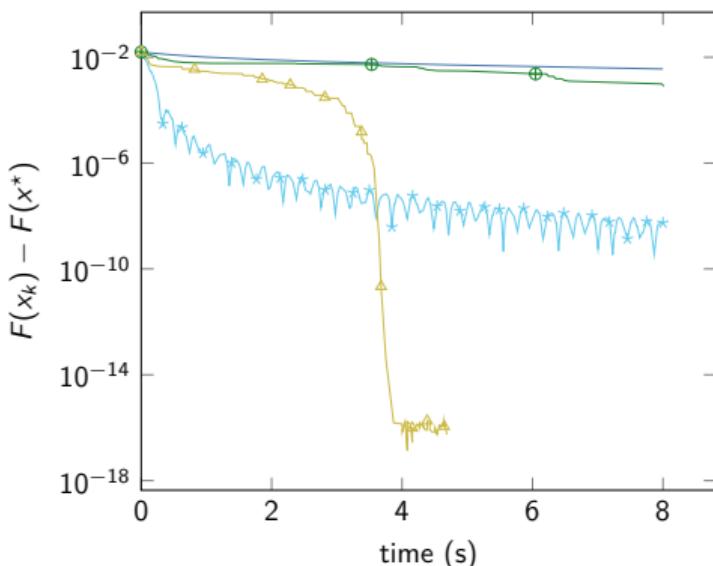
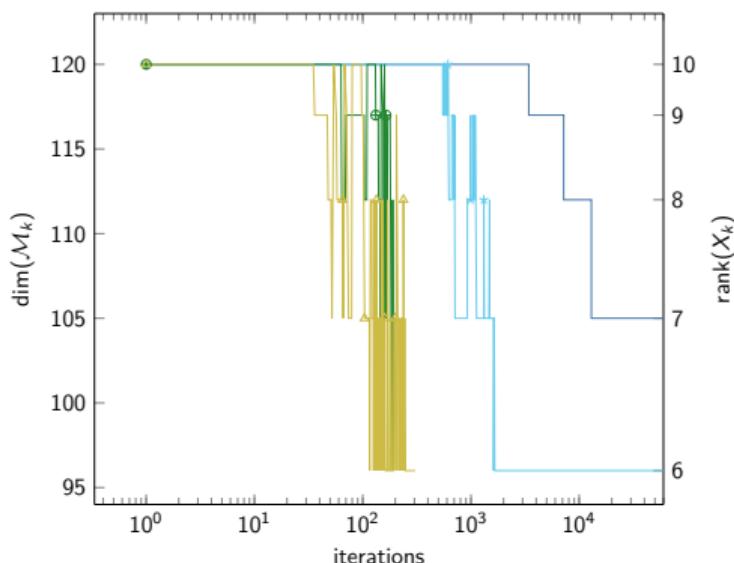
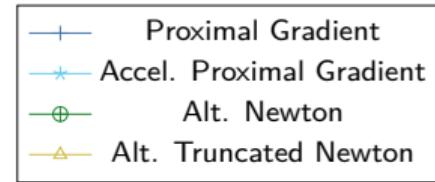
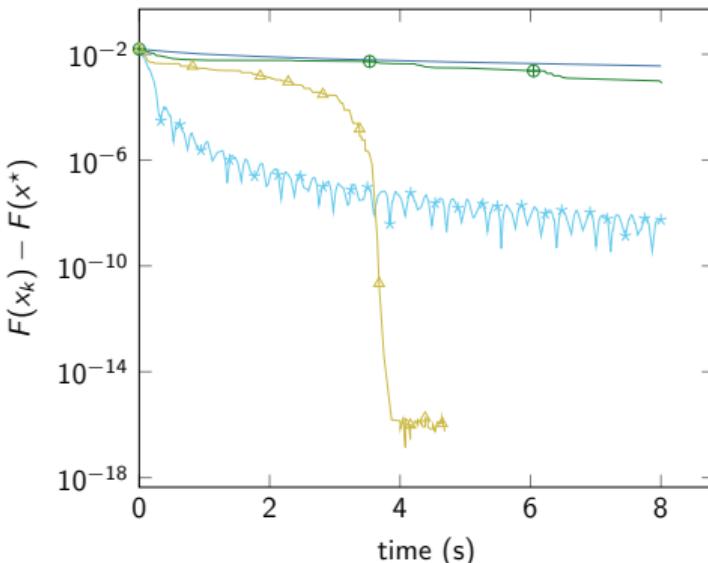


Illustration 2: tracenorm regression

$$\min_{X \in \mathbb{R}^{10 \times 12}} \sum_{i=1}^m (\langle A_i, X \rangle - y_i)^2 + \lambda \|X\|_*$$

In this instance, $m = 60$

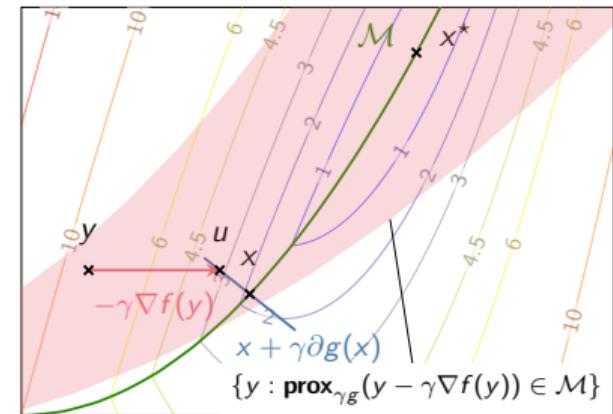
$$\mathcal{M}^* = \{X : \text{rank}(X) = 6\}$$



Summary on additive problems

Main messages:

- ▶ The proximal-gradient **identifies structure**
- ▶ With Riemannian Newton steps, we propose a global algorithm that
 - ▶ identifies \mathcal{M}^*
 - ▶ converges locally quadratically
- ▶ We observe these results numerically



Proximal gradient identification

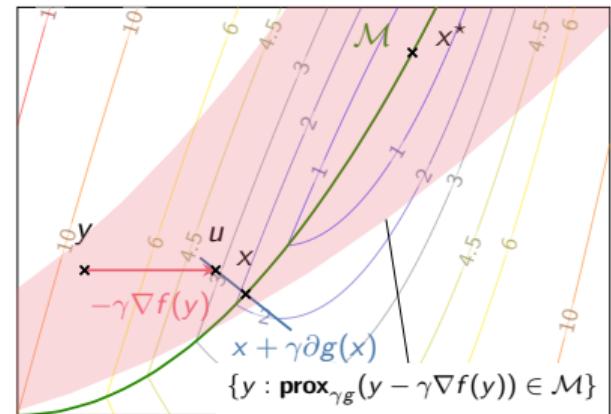
Summary on additive problems

Main messages:

- ▶ The proximal-gradient **identifies structure**
- ▶ With Riemannian Newton steps, we propose a global algorithm that
 - ▶ identifies \mathcal{M}^*
 - ▶ converges locally quadratically
- ▶ We observe these results numerically

Extensions: This is the simplest scheme

- ▶ other Riemannian method cubic regularization of Newton
- ▶ other interleaving



Proximal gradient identification

Introduction
○○○○○○○

Additive nonsmoothness $f + g$
○○○○○○○○○○

Composite nonsmoothness $g \circ c$
●○○○○○○○○○○○○○○

Conclusion
○○○○○

Outline

Introduction

Additive nonsmoothness $f + g$

Composite nonsmoothness $g \circ c$

Conclusion

Composite problem

Find $x^* \in \arg \min_{x \in \mathbb{R}^n} F(x) = g \circ c(x)$, with g nonsmooth and c a smooth mapping

Example:

- ▶ Maximum of smooth functions: $F(x) = \max(c_1(x), \dots, c_m(x))$ ◇ Womersley '86
- ▶ Spectral maximum: $F(x) = \lambda_{\max}(c(x))$, where $c(x) \in \mathbb{S}_m$ ◇ Oustry '00, Noll '05

Again, nonsmoothness is **explicit**: $\text{prox}_{\gamma g}$ is available

Composite problem

Find $x^* \in \arg \min_{x \in \mathbb{R}^n} F(x) = g \circ c(x)$, with g nonsmooth and c a smooth mapping

Example:

- ▶ Maximum of smooth functions: $F(x) = \max(c_1(x), \dots, c_m(x))$ ◇ Womersley '86
- ▶ Spectral maximum: $F(x) = \lambda_{\max}(c(x))$, where $c(x) \in \mathbb{S}_m$ ◇ Oustry '00, Noll '05

Again, nonsmoothness is **explicit**: $\text{prox}_{\gamma g}$ is available

- ▷ Structure detection: iterates of optimization methods **do not belong to manifolds** anymore
 - ▶ Composite bundle ◇ Sagastizábal '13
 - ▶ proxlinear methods ◇ Lewis Wright '16

→ What structure near point x ?

Back to the prox

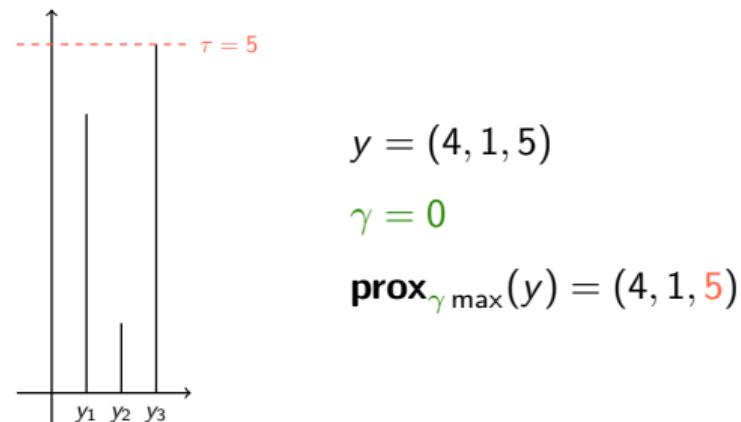
$$\text{prox}_{\gamma g}(y) = \arg \min_u \left\{ g(u) + \frac{1}{2\gamma} \|u - y\|^2 \right\}$$

Example: Prox of max function

$$[\text{prox}_{\gamma \max}(y)]_i = \begin{cases} \tau & \text{if } y_i \geq \tau \\ y_i & \text{else} \end{cases}$$

where τ solves $\sum_{\{i:y_i > \tau\}} (y_i - \tau) = \gamma$

$$\mathcal{M}_I = \{y : y_i = \max(y) \text{ for } i \in I\}$$



Back to the prox

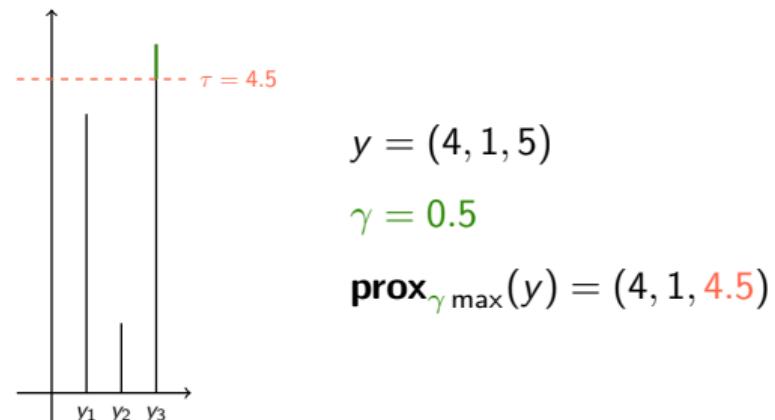
$$\text{prox}_{\gamma g}(y) = \arg \min_u \left\{ g(u) + \frac{1}{2\gamma} \|u - y\|^2 \right\}$$

Example: Prox of max function

$$[\text{prox}_{\gamma \max}(y)]_i = \begin{cases} \tau & \text{if } y_i \geq \tau \\ y_i & \text{else} \end{cases}$$

where τ solves $\sum_{\{i:y_i > \tau\}} (y_i - \tau) = \gamma$

$$\mathcal{M}_I = \{y : y_i = \max(y) \text{ for } i \in I\}$$



$$\gamma = 0.5$$

$$\text{prox}_{\gamma \max}(y) = (4, 1, 4.5)$$

Back to the prox

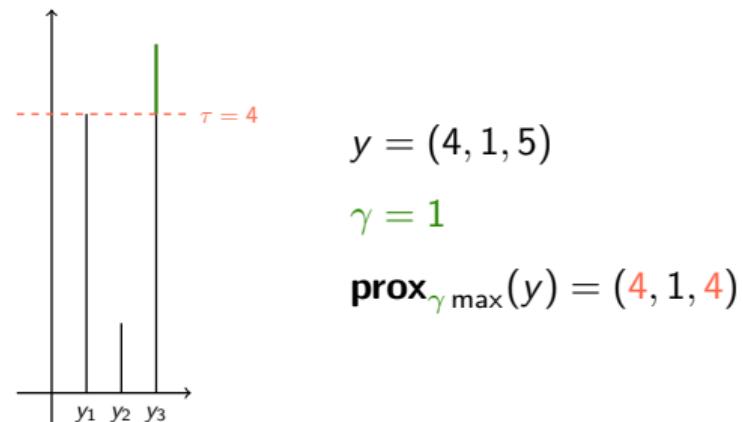
$$\text{prox}_{\gamma g}(y) = \arg \min_u \left\{ g(u) + \frac{1}{2\gamma} \|u - y\|^2 \right\}$$

Example: Prox of max function

$$[\text{prox}_{\gamma \max}(y)]_i = \begin{cases} \tau & \text{if } y_i \geq \tau \\ y_i & \text{else} \end{cases}$$

where τ solves $\sum_{\{i:y_i > \tau\}} (y_i - \tau) = \gamma$

$$\mathcal{M}_I = \{y : y_i = \max(y) \text{ for } i \in I\}$$



Back to the prox

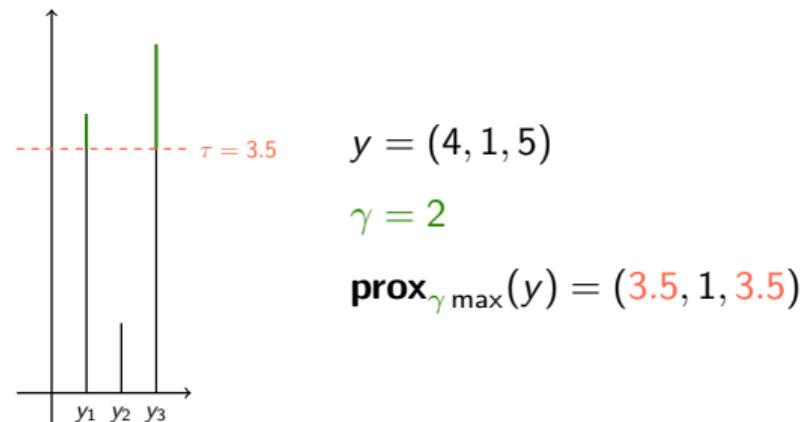
$$\text{prox}_{\gamma g}(y) = \arg \min_u \left\{ g(u) + \frac{1}{2\gamma} \|u - y\|^2 \right\}$$

Example: Prox of max function

$$[\text{prox}_{\gamma \max}(y)]_i = \begin{cases} \tau & \text{if } y_i \geq \tau \\ y_i & \text{else} \end{cases}$$

where τ solves $\sum_{\{i:y_i > \tau\}} (y_i - \tau) = \gamma$

$$\mathcal{M}_I = \{y : y_i = \max(y) \text{ for } i \in I\}$$



Back to the prox

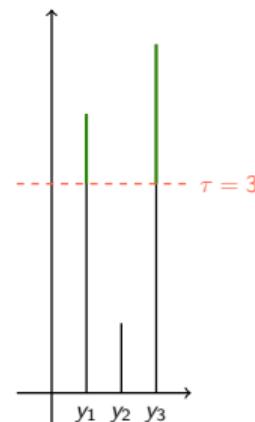
$$\text{prox}_{\gamma g}(y) = \arg \min_u \left\{ g(u) + \frac{1}{2\gamma} \|u - y\|^2 \right\}$$

Example: Prox of max function

$$[\text{prox}_{\gamma \max}(y)]_i = \begin{cases} \tau & \text{if } y_i \geq \tau \\ y_i & \text{else} \end{cases}$$

where τ solves $\sum_{\{i:y_i > \tau\}} (y_i - \tau) = \gamma$

$$\mathcal{M}_I = \{y : y_i = \max(y) \text{ for } i \in I\}$$



$$y = (4, 1, 5)$$

$$\gamma = 3$$

$$\text{prox}_{\gamma \max}(y) = (3, 1, 3)$$

Back to the prox

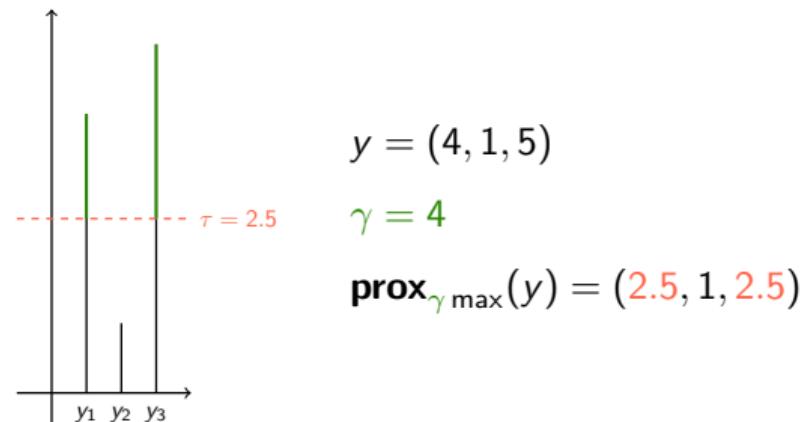
$$\text{prox}_{\gamma g}(y) = \arg \min_u \left\{ g(u) + \frac{1}{2\gamma} \|u - y\|^2 \right\}$$

Example: Prox of max function

$$[\text{prox}_{\gamma \max}(y)]_i = \begin{cases} \tau & \text{if } y_i \geq \tau \\ y_i & \text{else} \end{cases}$$

where τ solves $\sum_{\{i:y_i > \tau\}} (y_i - \tau) = \gamma$

$$\mathcal{M}_I = \{y : y_i = \max(y) \text{ for } i \in I\}$$



Back to the prox

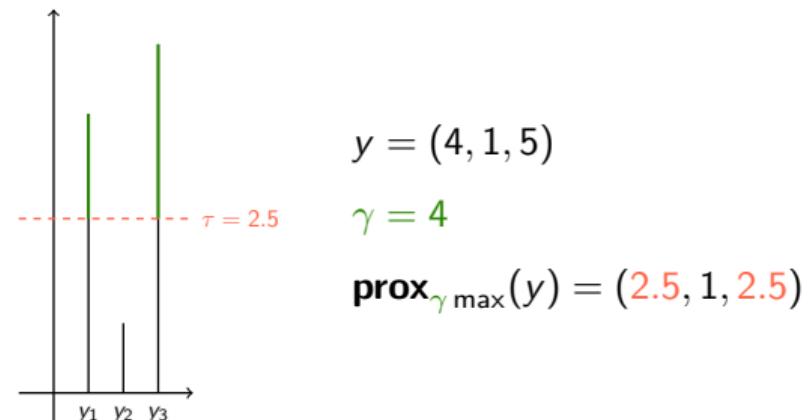
$$\text{prox}_{\gamma g}(y) = \arg \min_u \left\{ g(u) + \frac{1}{2\gamma} \|u - y\|^2 \right\}$$

Example: Prox of max function

$$[\text{prox}_{\gamma \max}(y)]_i = \begin{cases} \tau & \text{if } y_i \geq \tau \\ y_i & \text{else} \end{cases}$$

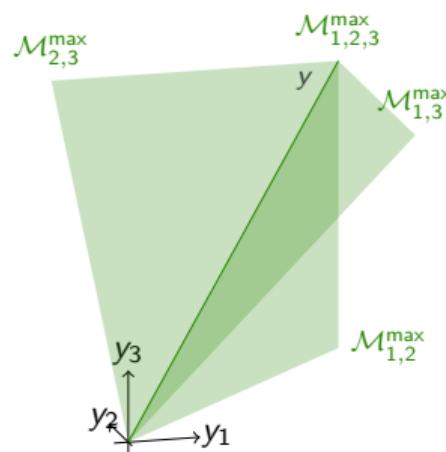
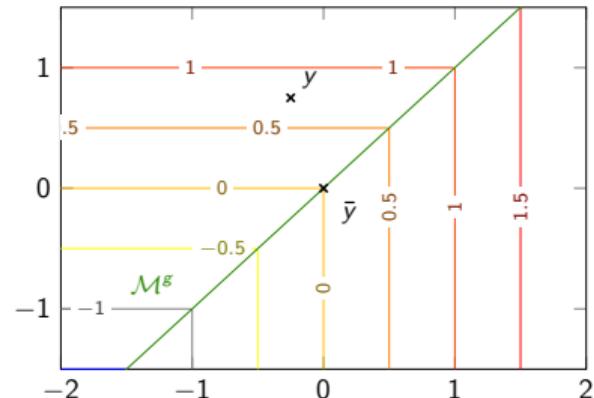
where τ solves $\sum_{\{i:y_i > \tau\}} (y_i - \tau) = \gamma$

$$\mathcal{M}_I = \{y : y_i = \max(y) \text{ for } i \in I\}$$

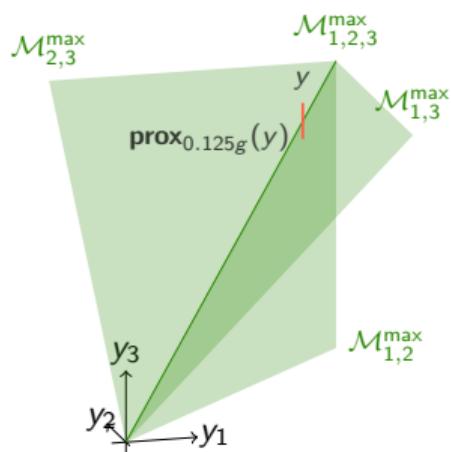
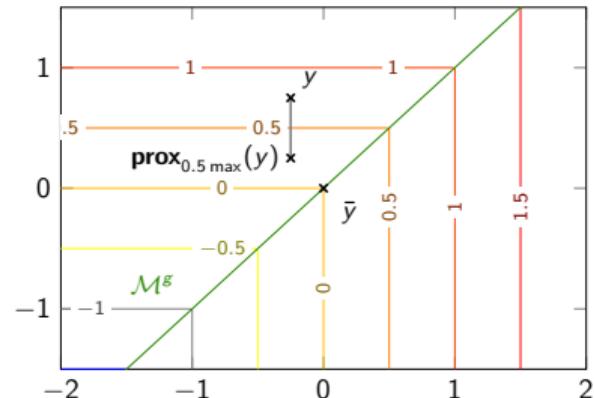


→ Computing $\text{prox}_{\gamma g}(y)$ also gives **structure information** $\mathcal{M} \ni \text{prox}_{\gamma g}(y)$.

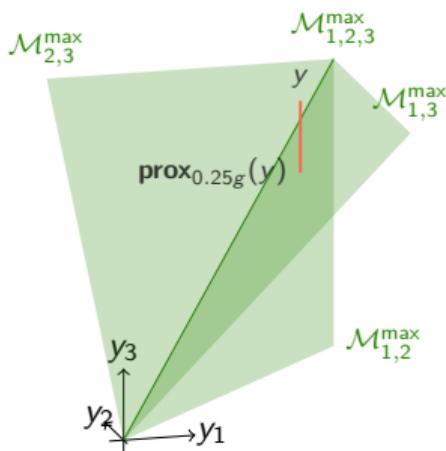
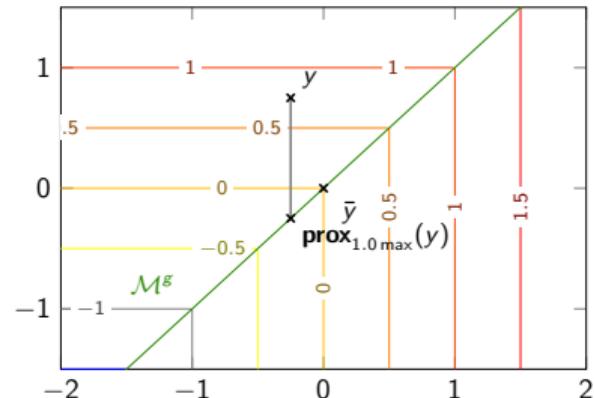
Identification for g , explicit $\text{prox}_{\gamma g}$



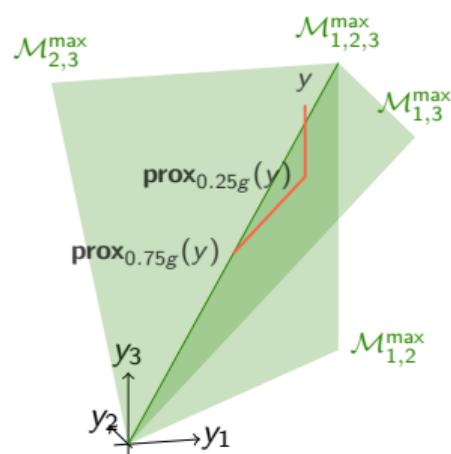
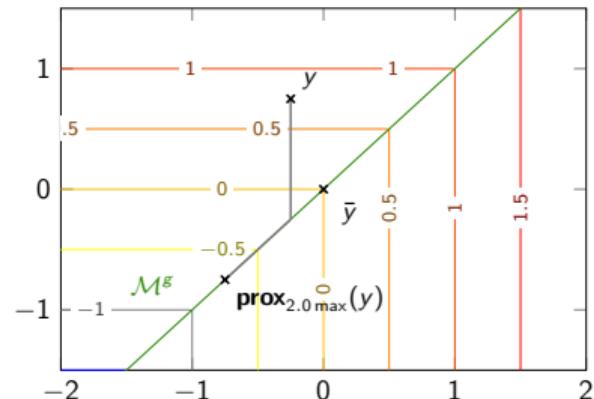
Identification for g , explicit $\text{prox}_{\gamma g}$



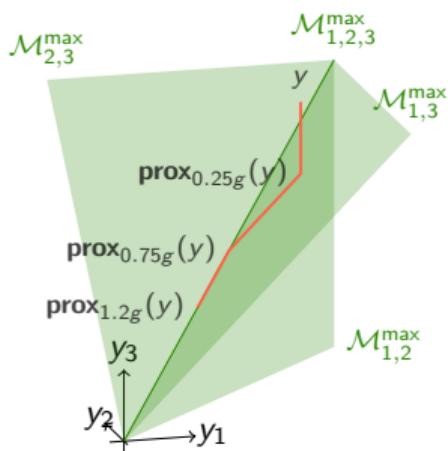
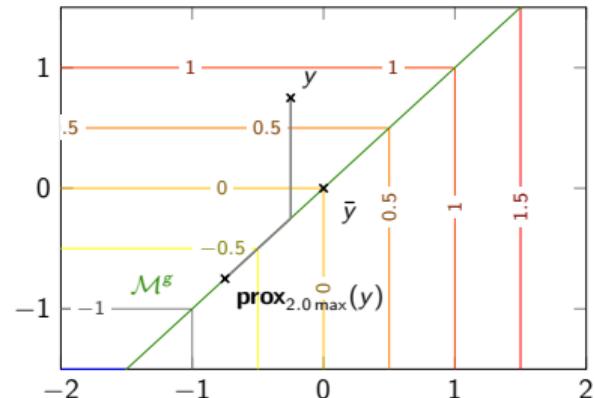
Identification for g , explicit $\text{prox}_{\gamma g}$



Identification for g , explicit $\text{prox}_{\gamma g}$



Identification for g , explicit $\text{prox}_{\gamma g}$



Identification for g , explicit $\text{prox}_{\gamma g}$

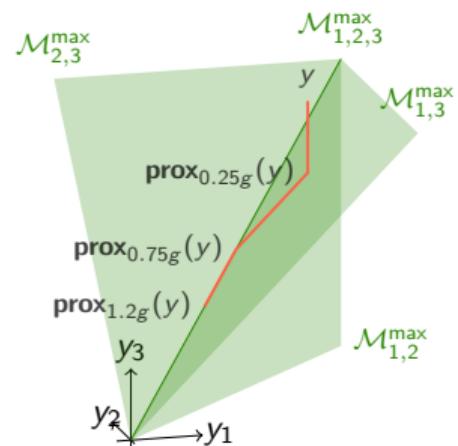
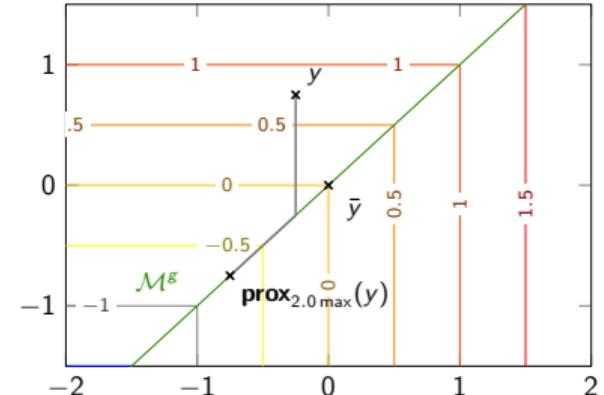
Lemma

Consider a function g and point \bar{y} with structure \mathcal{M}^g that meet two technical assumptions. For all y near \bar{y} ,

$$\text{prox}_{\gamma g}(y) \in \mathcal{M}^g \quad \text{for all } \gamma \in [\varphi^g(\text{dist}_{\mathcal{M}^g}(y)), \Gamma^g]$$

where $\Gamma^g > 0$ and $\varphi^g(t) = \frac{1}{c_{ri}}t + \mathcal{O}(t^2)$.

Technical assumptions: normal ascent, control on projection curves



Identification for g , explicit $\text{prox}_{\gamma g}$

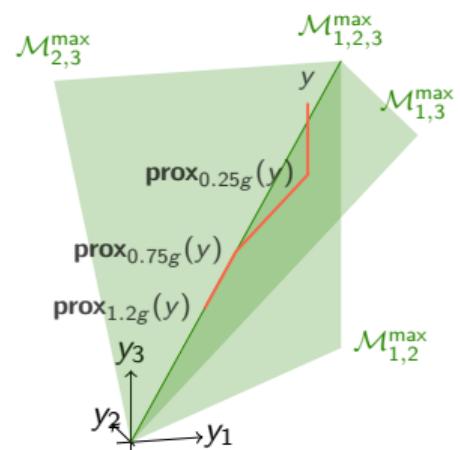
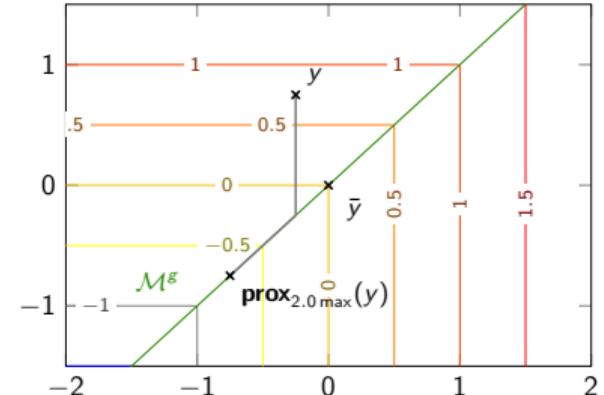
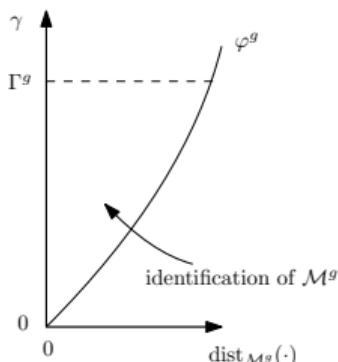
Lemma

Consider a function g and point \bar{y} with structure \mathcal{M}^g that meet two technical assumptions. For all y near \bar{y} ,

$$\text{prox}_{\gamma g}(y) \in \mathcal{M}^g \quad \text{for all } \gamma \in [\varphi^g(\text{dist}_{\mathcal{M}^g}(y)), \Gamma^g]$$

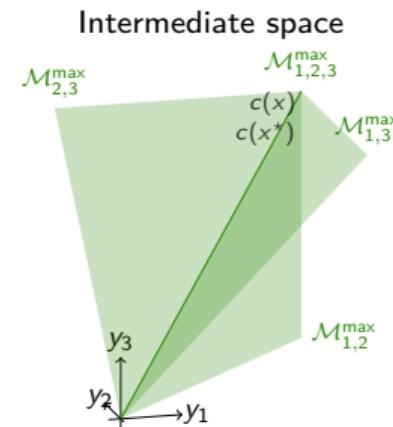
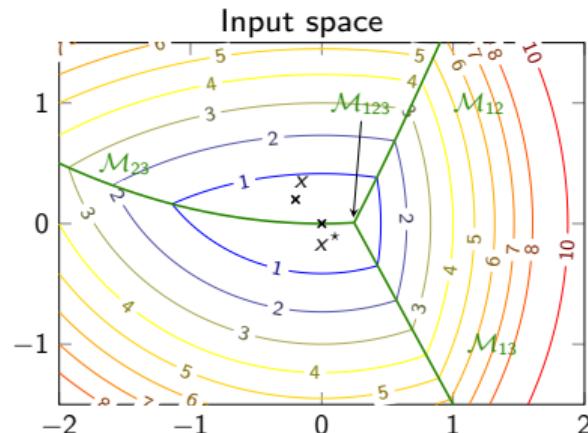
where $\Gamma^g > 0$ and $\varphi^g(t) = \frac{1}{c_{ri}}t + \mathcal{O}(t^2)$.

Technical assumptions: normal ascent, control on projection curves



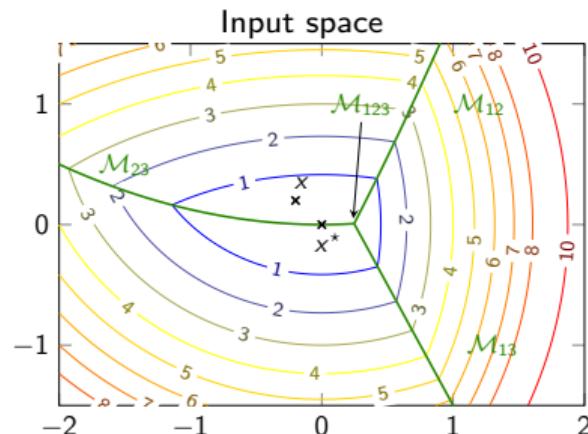
Identification for $g \circ c$, no prox $_{\gamma g \circ c}$

The prox of $F = g \circ c$ is **not available**, but we do have prox $_{\gamma g}$.

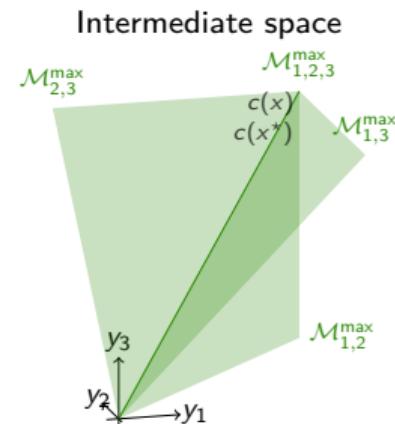


Identification for $g \circ c$, no prox _{$\gamma g \circ c$}

The prox of $F = g \circ c$ is **not available**, but we do have prox _{γg} .



$$F(x) = \max(c_1(x), c_2(x), c_3(x))$$



$$g(y) = \max(y_1, y_2, y_3)$$

Observation: prox _{γg} can map points to M^g .

The structure naturally lies in the intermediate space ... but that's ok!

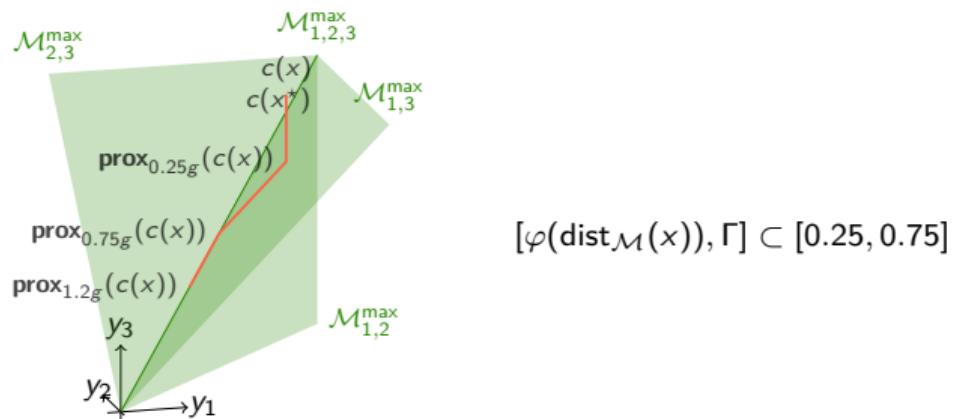
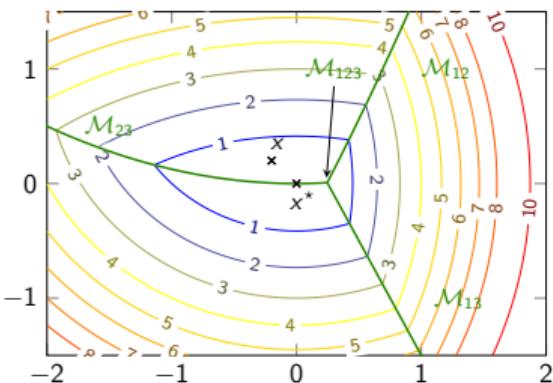
Identification for $g \circ c$

Theorem

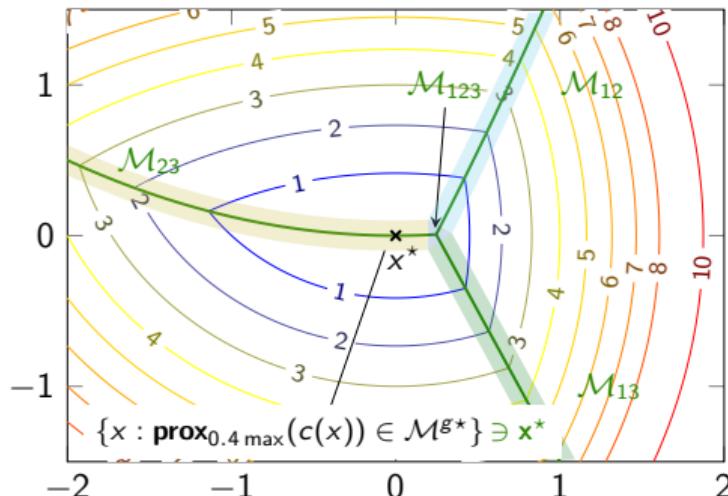
Consider g, c smooth and a point \bar{x} such that $c(\bar{x})$ has structure manifold \mathcal{M}^g with c and \mathcal{M}^g are transversal at $c(\bar{x})$. For all x near \bar{x} ,

$$\text{prox}_{\gamma g}(c(x)) \in \mathcal{M}^g \quad \text{for all } \gamma \in [\varphi(\text{dist}_{\mathcal{M}}(x)), \Gamma]$$

where $\Gamma > 0$ and $\varphi(t) = \frac{c_{\text{map}}}{c_{ri}}t + \mathcal{O}(t^2)$. Furthermore, $\mathcal{M} = c^{-1}(\mathcal{M}^g)$.



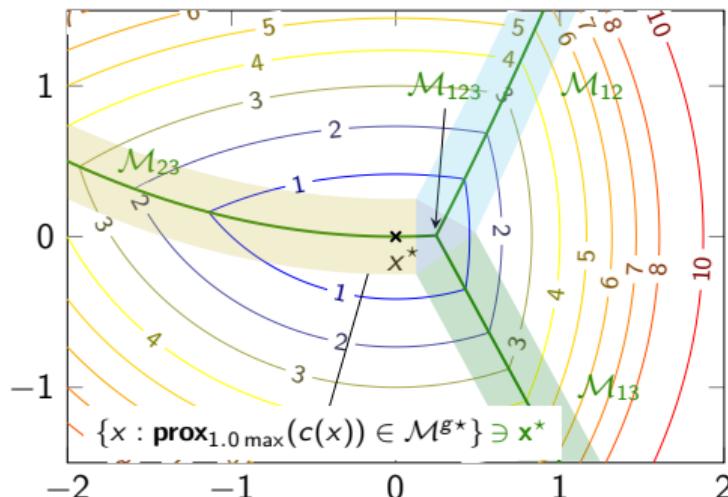
Towards an algorithm: i. detecting structure



γ too small \Rightarrow detection of \mathcal{M}^* only near x^* ;

$$\mathcal{M}^* = \mathcal{M}_{23}$$

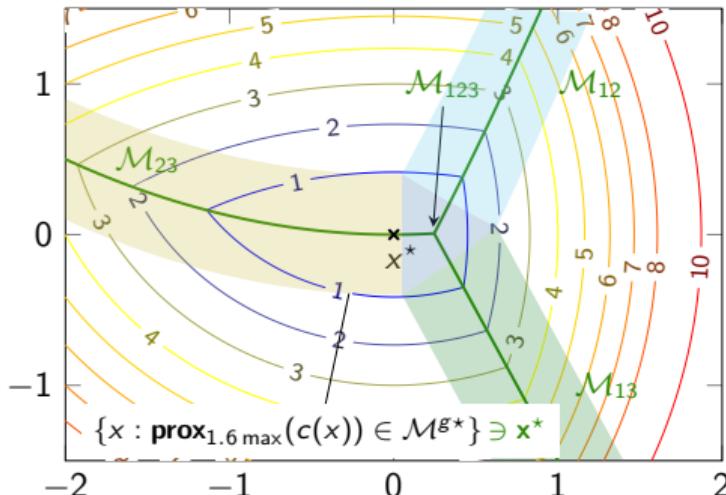
Towards an algorithm: i. detecting structure



γ too small \Rightarrow detection of \mathcal{M}^* only near x^* ;

$$\mathcal{M}^* = \mathcal{M}_{23}$$

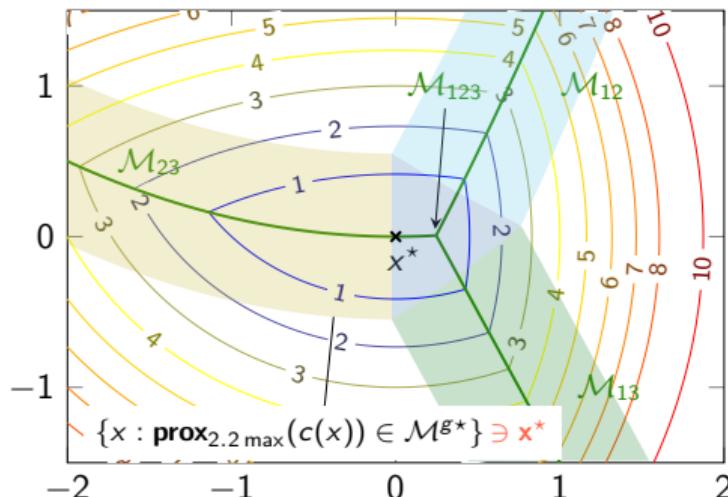
Towards an algorithm: i. detecting structure



γ too small \Rightarrow detection of \mathcal{M}^* only near x^* ;

$$\mathcal{M}^* = \mathcal{M}_{23}$$

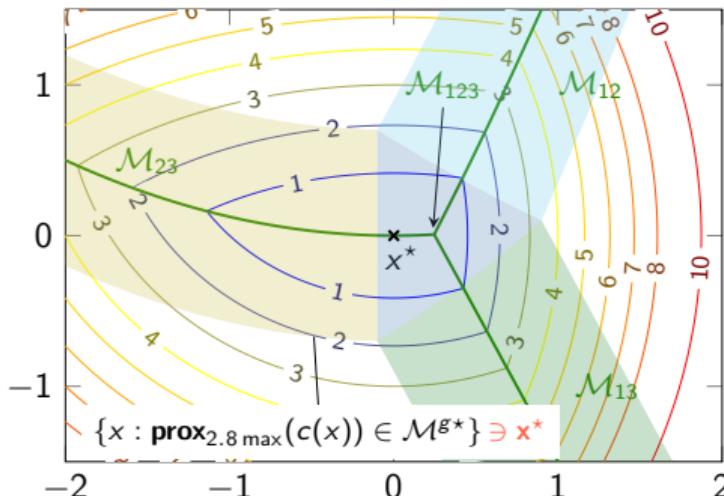
Towards an algorithm: i. detecting structure



γ too small \Rightarrow detection of \mathcal{M}^* only near x^* ;

$$\mathcal{M}^* = \mathcal{M}_{23}$$

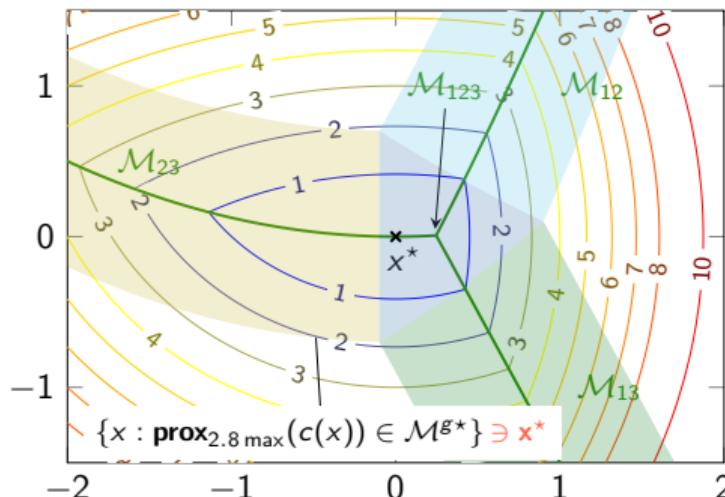
Towards an algorithm: i. detecting structure



γ too small \Rightarrow detection of \mathcal{M}^* only near x^* ;

$$\mathcal{M}^* = \mathcal{M}_{23}$$

Towards an algorithm: i. detecting structure

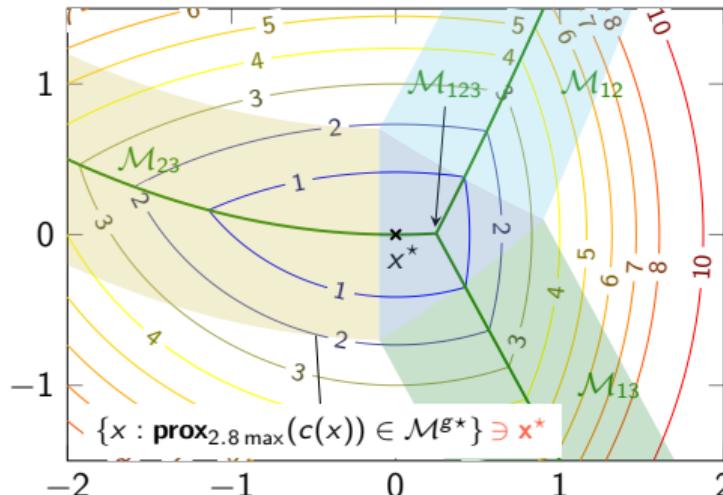


$$\mathcal{M}^* = \mathcal{M}_{23}$$

γ too small \Rightarrow detection of \mathcal{M}^* only near x^* ;

γ too large \Rightarrow no detection of \mathcal{M}^* near x^* .

Towards an algorithm: i. detecting structure



$$\mathcal{M}^* = \mathcal{M}_{23}$$

Bottom line: $\text{prox}_{\gamma g} \circ c(\cdot)$ detects \mathcal{M}^* from any x near x^* when $\gamma \in [\frac{c_{map}}{c_{ri}} \text{dist}_{\mathcal{M}}(x), \Gamma]$.

→ How to choose the step in practice?

Towards an algorithm: ii. exploiting structure

\mathcal{M} is smooth $\exists h$ smooth s.t. $x \in \mathcal{M} \Leftrightarrow h(x) = 0$
 F smooth on \mathcal{M} $\exists \tilde{F}$ smooth s.t. $F|_{\mathcal{M}} \equiv \tilde{F}$ on \mathcal{M}

$$\min_{x \in \mathbb{R}^n} F(x) \quad \xrightarrow{\text{with } \mathcal{M}} \quad \min_x \tilde{F}(x) \text{ s.t. } h(x) = 0.$$

Towards an algorithm: ii. exploiting structure

$$\begin{array}{lll} \mathcal{M} \text{ is smooth} & \exists h \text{ smooth s.t. } x \in \mathcal{M} \Leftrightarrow h(x) = 0 \\ F \text{ smooth on } \mathcal{M} & \exists \tilde{F} \text{ smooth s.t. } F|_{\mathcal{M}} \equiv \tilde{F} \text{ on } \mathcal{M} \end{array}$$

$$\min_{x \in \mathbb{R}^n} F(x) \quad \xrightarrow{\text{with } \mathcal{M}} \quad \min_x \tilde{F}(x) \text{ s.t. } h(x) = 0.$$

Example ($F = \max(c_1, c_2, c_3)$)

For structure \mathcal{M}_{12} ,

- ▶ $h = c_1 - c_2$
- ▶ $\tilde{F}(x) = (c_1 + c_2)/2$

Towards an algorithm: ii. exploiting structure

$$\begin{array}{lll} \mathcal{M} \text{ is smooth} & \exists h \text{ smooth s.t. } x \in \mathcal{M} \Leftrightarrow h(x) = 0 \\ F \text{ smooth on } \mathcal{M} & \exists \tilde{F} \text{ smooth s.t. } F|_{\mathcal{M}} \equiv \tilde{F} \text{ on } \mathcal{M} \end{array}$$

$$\min_{x \in \mathbb{R}^n} F(x) \quad \xrightarrow{\text{with } \mathcal{M}} \quad \min_x \tilde{F}(x) \text{ s.t. } h(x) = 0.$$

Example ($F = \max(c_1, c_2, c_3)$)

For structure \mathcal{M}_{12} ,

- ▶ $h = c_1 - c_2$
- ▶ $\tilde{F}(x) = (c_1 + c_2)/2$

▷ Many tools for smooth constrained optimization: Interior Point Methods, **Sequential Quadratic Programming**, Augmented Lagrangian Methods, etc

Algorithm: detection & exploitation

Iteration k :

Algorithm: detection & exploitation

Iteration k :

- ▷ Compute $\text{prox}_{\gamma_k g}(c(x_k))$ and obtain \mathcal{M}_k and h_k, \tilde{F}_k

Algorithm: detection & exploitation

Iteration k :

▷ Compute $\text{prox}_{\gamma_k g}(c(x_k))$ and obtain \mathcal{M}_k and h_k, \tilde{F}_k

▷ **Newton-SQP step** relative to \mathcal{M}_k :

$$\begin{aligned} d_k^{\text{SQP}}(x_k) &= \arg \min_{d \in \mathbb{R}^n} \quad \langle \nabla \tilde{F}_k(x_k), d \rangle + \frac{1}{2} \langle \nabla_{xx}^2 L_k(x_k, \lambda_k(x_k)) d, d \rangle \\ \text{s.t.} \quad h_k(x_k) + D h_k(x_k) d &= 0 \end{aligned}$$

where $L_k(x, \lambda) = \tilde{F}_k(x) + \langle \lambda, h_k(x) \rangle$, and $\lambda_k(x_k) = \arg \min_{\lambda \in \mathbb{R}^r} \left\| \nabla \tilde{F}_k(x_k) + \sum_{i=1}^m \lambda_i \nabla h_{k,i}(x_k) \right\|^2$

Set $x_{k+1} = x_k + d_k^{\text{SQP}}(x_k)$ if $F(x_k + d_k^{\text{SQP}}(x_k)) < F(x_k)$

Algorithm: detection & exploitation

Iteration k :

▷ Compute $\text{prox}_{\gamma_k g}(c(x_k))$ and obtain \mathcal{M}_k and h_k, \tilde{F}_k

▷ **Newton-SQP step** relative to \mathcal{M}_k :

$$\begin{aligned} d_k^{\text{SQP}}(x_k) &= \arg \min_{d \in \mathbb{R}^n} \quad \langle \nabla \tilde{F}_k(x_k), d \rangle + \frac{1}{2} \langle \nabla_{xx}^2 L_k(x_k, \lambda_k(x_k)) d, d \rangle \\ \text{s.t.} \quad h_k(x_k) + D h_k(x_k) d &= 0 \end{aligned}$$

where $L_k(x, \lambda) = \tilde{F}_k(x) + \langle \lambda, h_k(x) \rangle$, and $\lambda_k(x_k) = \arg \min_{\lambda \in \mathbb{R}^r} \left\| \nabla \tilde{F}_k(x_k) + \sum_{i=1}^m \lambda_i \nabla h_{k,i}(x_k) \right\|^2$

Set $x_{k+1} = x_k + d_k^{\text{SQP}}(x_k)$ if $F(x_k + d_k^{\text{SQP}}(x_k)) < F(x_k)$

▷ $\gamma_{k+1} = \frac{\gamma_k}{2}$

Local exact structure identification and quadratic convergence

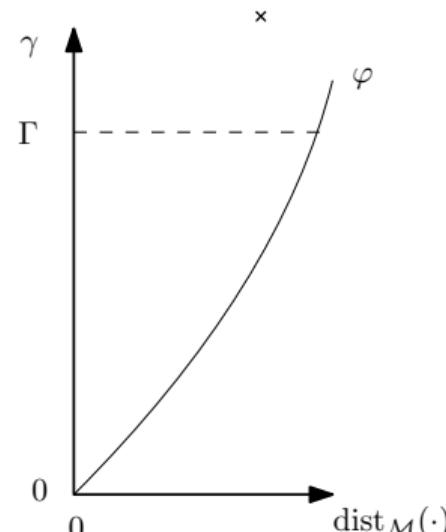
Theorem

Consider $F = g \circ c$ and a minimizer $x^* \in \mathcal{M}^*$ that meet some assumptions

Then, if x_0 is close enough to x^* and γ_0 is large enough, after some finite time

- ▶ $\mathcal{M}_k = \mathcal{M}^*$

- ▶ x_k converges to x^* at a **quadratic rate**: $\|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2$



Sketch of proof, on a picture:

Local exact structure identification and quadratic convergence

Theorem

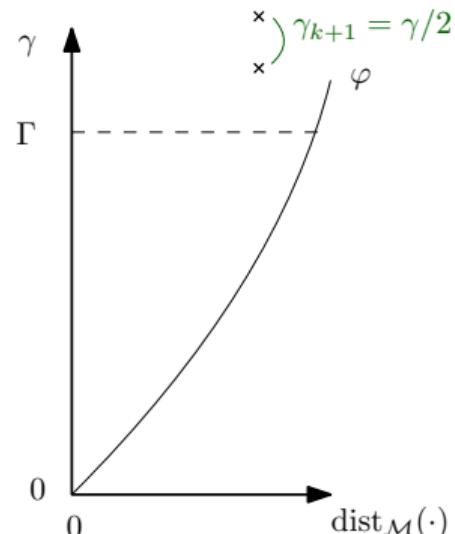
Consider $F = g \circ c$ and a minimizer $x^* \in \mathcal{M}^*$ that meet some assumptions

Then, if x_0 is close enough to x^* and γ_0 is large enough, after some finite time

► $\mathcal{M}_k = \mathcal{M}^*$

► x_k converges to x^* at a **quadratic rate**: $\|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2$

Sketch of proof, on a picture:



Local exact structure identification and quadratic convergence

Theorem

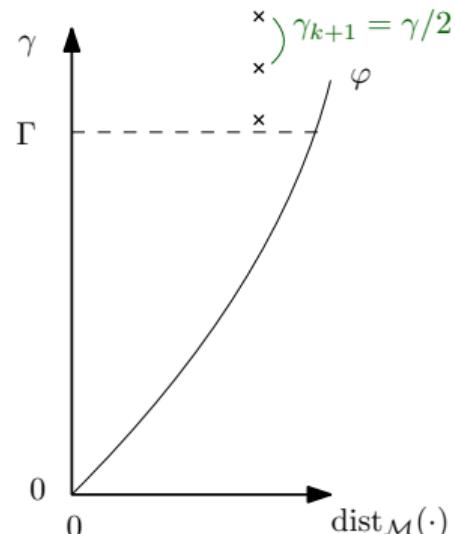
Consider $F = g \circ c$ and a minimizer $x^* \in \mathcal{M}^*$ that meet some assumptions

Then, if x_0 is close enough to x^* and γ_0 is large enough, after some finite time

- ▶ $\mathcal{M}_k = \mathcal{M}^*$

- ▶ x_k converges to x^* at a **quadratic rate**: $\|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2$

Sketch of proof, on a picture:



Local exact structure identification and quadratic convergence

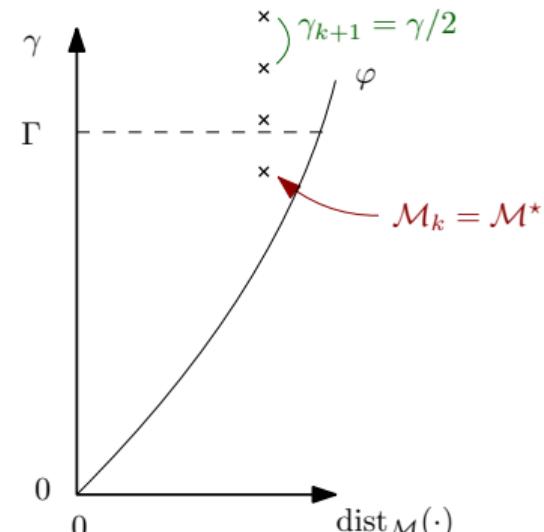
Theorem

Consider $F = g \circ c$ and a minimizer $x^* \in \mathcal{M}^*$ that meet some assumptions

Then, if x_0 is close enough to x^* and γ_0 is large enough, after some finite time

► $\mathcal{M}_k = \mathcal{M}^*$

► x_k converges to x^* at a **quadratic rate**: $\|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2$



Sketch of proof, on a picture:

Local exact structure identification and quadratic convergence

Theorem

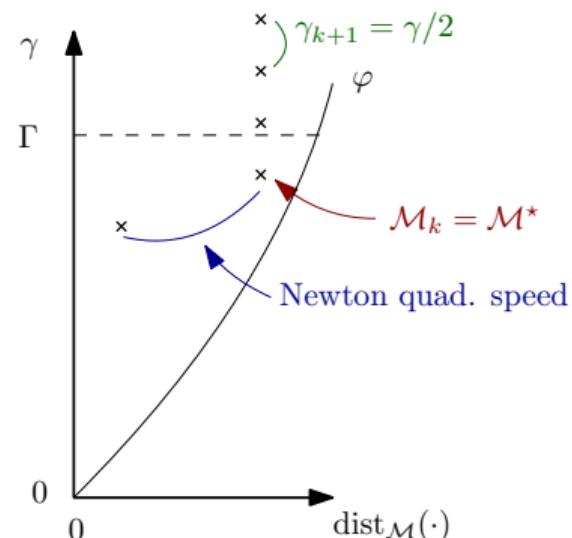
Consider $F = g \circ c$ and a minimizer $x^* \in \mathcal{M}^*$ that meet some assumptions

Then, if x_0 is close enough to x^* and γ_0 is large enough, after some finite time

► $\mathcal{M}_k = \mathcal{M}^*$

► x_k converges to x^* at a **quadratic rate**: $\|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2$

Sketch of proof, on a picture:



Local exact structure identification and quadratic convergence

Theorem

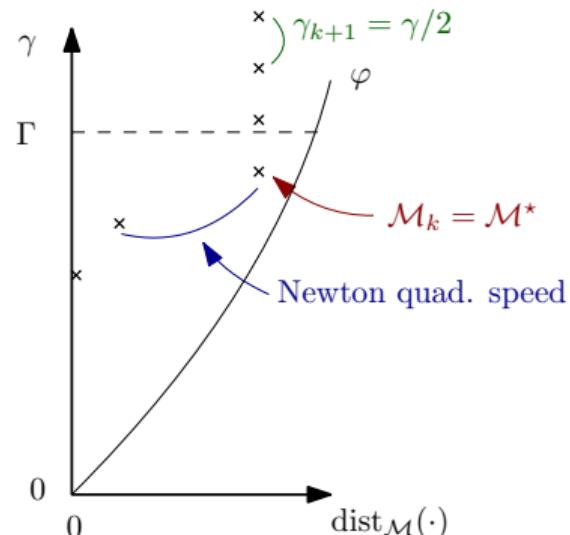
Consider $F = g \circ c$ and a minimizer $x^* \in \mathcal{M}^*$ that meet some assumptions

Then, if x_0 is close enough to x^* and γ_0 is large enough, after some finite time

► $\mathcal{M}_k = \mathcal{M}^*$

► x_k converges to x^* at a **quadratic rate**: $\|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2$

Sketch of proof, on a picture:



Local exact structure identification and quadratic convergence

Theorem

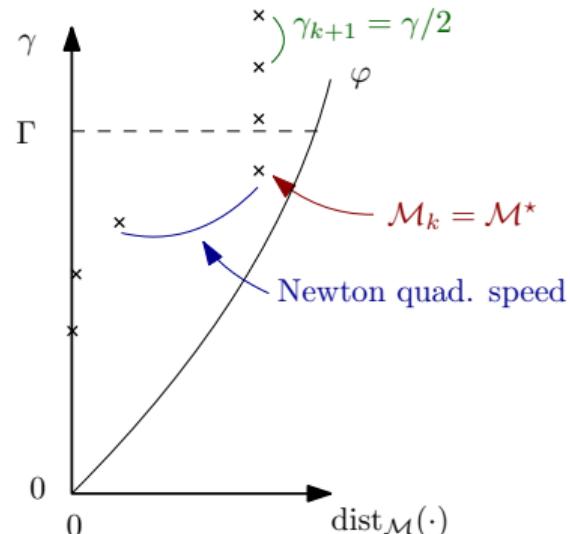
Consider $F = g \circ c$ and a minimizer $x^* \in \mathcal{M}^*$ that meet some assumptions

Then, if x_0 is close enough to x^* and γ_0 is large enough, after some finite time

► $\mathcal{M}_k = \mathcal{M}^*$

► x_k converges to x^* at a **quadratic rate**: $\|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2$

Sketch of proof, on a picture:

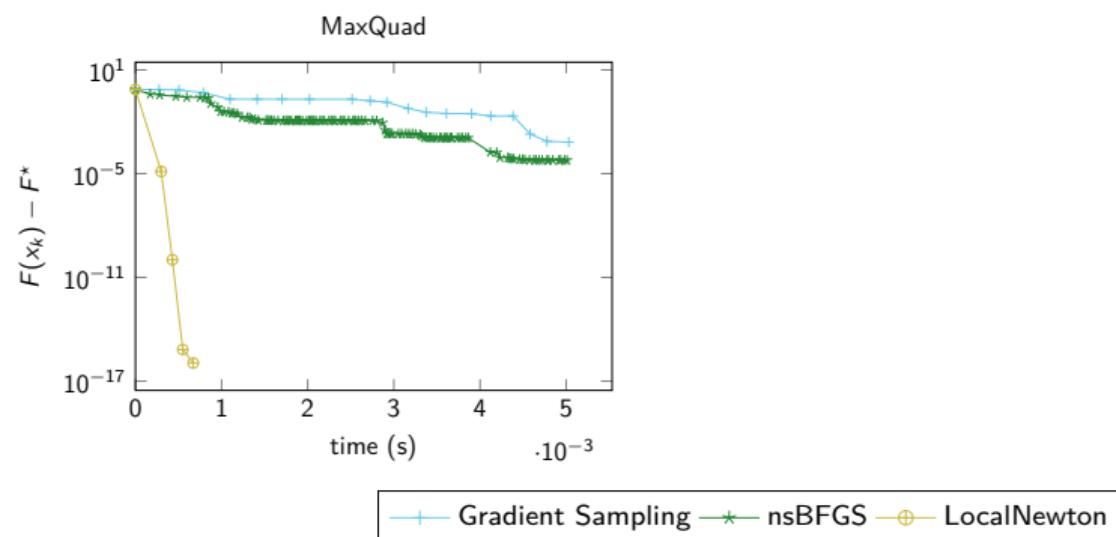


Illustrations: local behavior

$$\min_{x \in \mathbb{R}^{10}} \max_{i=1, \dots, 5} (c_i(x))$$

In this historical instance ◇ HULL '93

$$\mathcal{M}^* = \{x : c_2(x) = \dots = c_5(x)\}$$



Illustrations: local behavior

$$\min_{x \in \mathbb{R}^{10}} \max_{i=1, \dots, 5} (c_i(x))$$

In this historical instance ◇ HULL '93

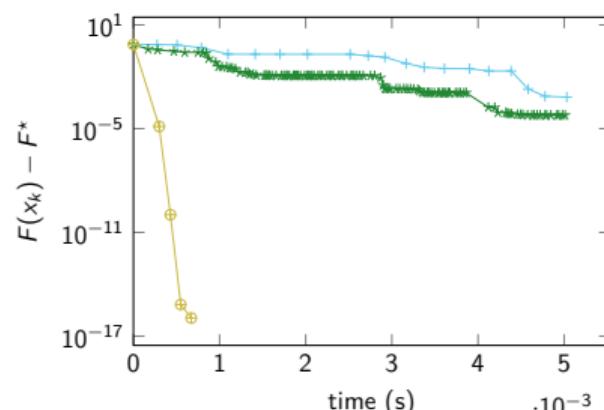
$$\mathcal{M}^* = \{x : c_2(x) = \dots = c_5(x)\}$$

$$\min_{x \in \mathbb{R}^n} \lambda_{\max} \left(A_0 + \sum_{i=1}^n x_i A_i \right)$$

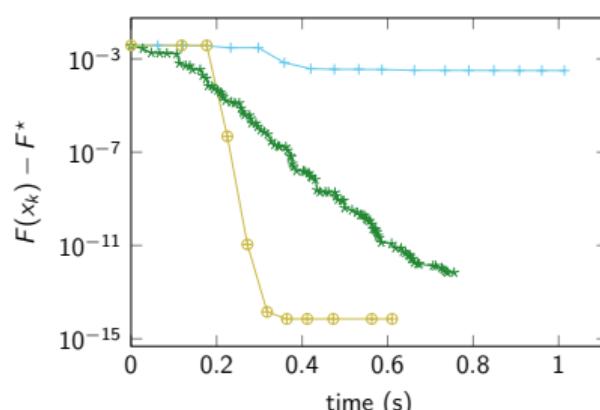
In this instance, $n = 25$, $A_i \in \mathbb{S}_{50}$

$$\mathcal{M}^* = \{x : \lambda_{\max}(c(x)) \text{ has multiplicity 3}\}$$

MaxQuad



Eigmax



Gradient Sampling — nsBFGS — LocalNewton

So far, on composite problems

Main messages:

- ▶ The proximal operator **identifies structure** in composite problems
- ▶ With Newton SQP steps, we propose a algorithm that **locally**
 - ▶ identifies \mathcal{M}^*
 - ▶ converges quadratically
- ▶ We observe these results numerically

So far, on composite problems

Main messages:

- ▶ The proximal operator **identifies structure** in composite problems
- ▶ With Newton SQP steps, we propose a algorithm that **locally**
 - ▶ identifies \mathcal{M}^*
 - ▶ converges quadratically
- ▶ We observe these results numerically

So far, on composite problems

Main messages:

- ▶ The proximal operator **identifies structure** in composite problems
- ▶ With Newton SQP steps, we propose a algorithm that **locally**
 - ▶ identifies \mathcal{M}^*
 - ▶ converges quadratically
- ▶ We observe these results numerically

Do we really need to start near x^* ?

Towards a global nonsmooth Newton method

- ▷ **Difficulty # 1:** Newton's direction may fail to provide functional descent.
→ linesearch: find α such that $F(x + \alpha d) \leq F(x) + \alpha m F'(x; d)$.
 - ▶ Is d a descent direction $F'(x; d) < 0$?
 - ▶ The linesearch may jeopardizes the final quadratic rate? Maratos effect, etc

Towards a global nonsmooth Newton method

▷ **Difficulty # 1:** Newton's direction may fail to provide functional descent.

→ linesearch: find α such that $F(x + \alpha d) \leq F(x) + \alpha m F'(x; d)$.

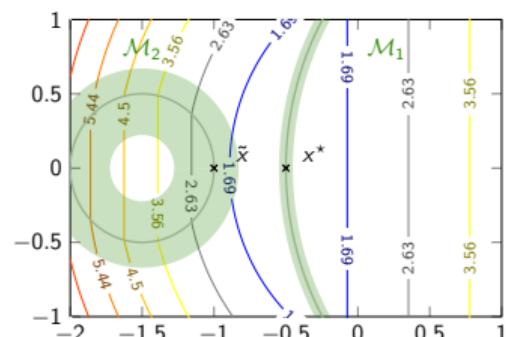
- ▶ Is d a descent direction $F'(x; d) < 0$?
- ▶ The linesearch may jeopardizes the final quadratic rate? Maratos effect, etc

▷ **Difficulty # 2:** Detection of relevant structure, far from x^*

Example In green, structure detection of $\text{prox}_{\gamma g} \circ c$

Point \tilde{x}

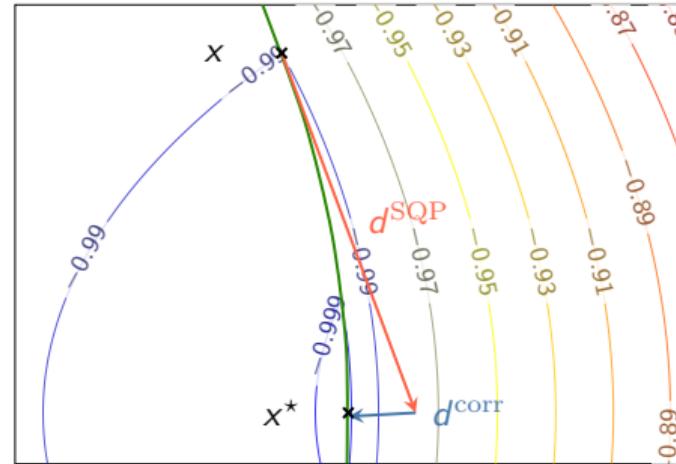
- ▶ is not a minimizer
- ▶ but stable structure detection by $\text{prox}_{\gamma g} \circ c$.



No Maratos effect with correction

Maratos effect: $F(x + d^{\text{SQP}}) > F(x)$
 \rightarrow Newton direction is not accepted

Second-order correction: from $x + d^{\text{SQP}}$, go towards \mathcal{M} Newton-Raphson step



Theorem (Local admissibility of unit stepsize)

Consider $F = g \circ c$ and a strong minimizer $x^* \in \mathcal{M}^*$ that meet some assumptions
Then if x is near x^* and $m \in (0, 1/2)$,

$$F(x + d^{\text{SQP}} + d^{\text{corr}}) \leq F(x) + mF'(x; d^{\text{SQP}}).$$

Heuristic algorithm illustrated

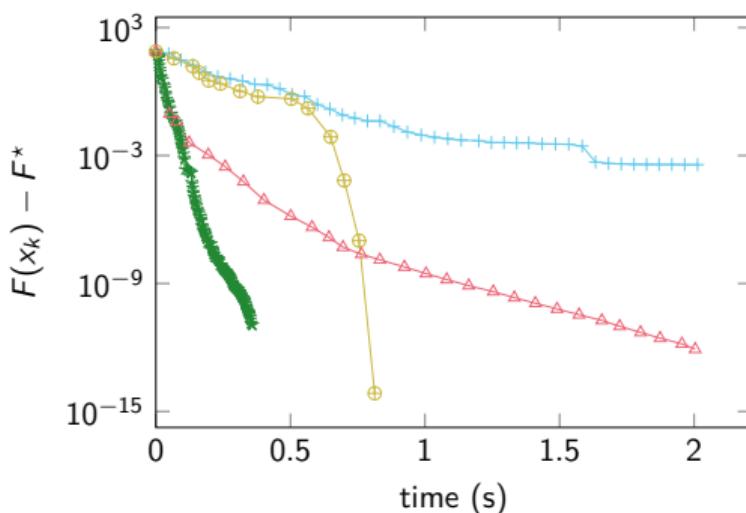
$$\min_{x \in \mathbb{R}^n} \lambda_{\max} \left(A_0 + \sum_{i=1}^n x_i A_i \right)$$

In this instance, $n = 25$, $A_i \in \mathbb{S}_{50}$

$$\mathcal{M}^* = \{x : \lambda_{\max}(c(x)) \text{ has multiplicity 3}\}$$

- +— Gradient Sampling
- *— nsBFGS
- ⊕— heuristic global Newton
- △— \mathcal{VU} -Newton bundle

(our implementation of ◊ Mifflin Sagastizábal '05)

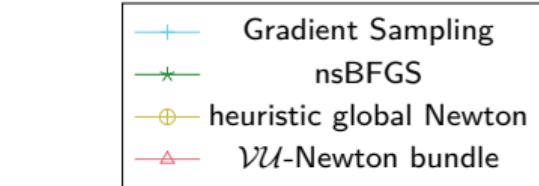
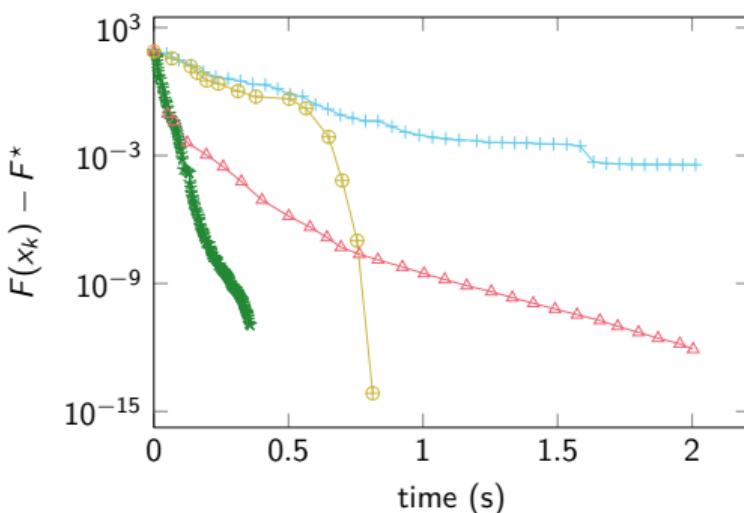


Heuristic algorithm illustrated

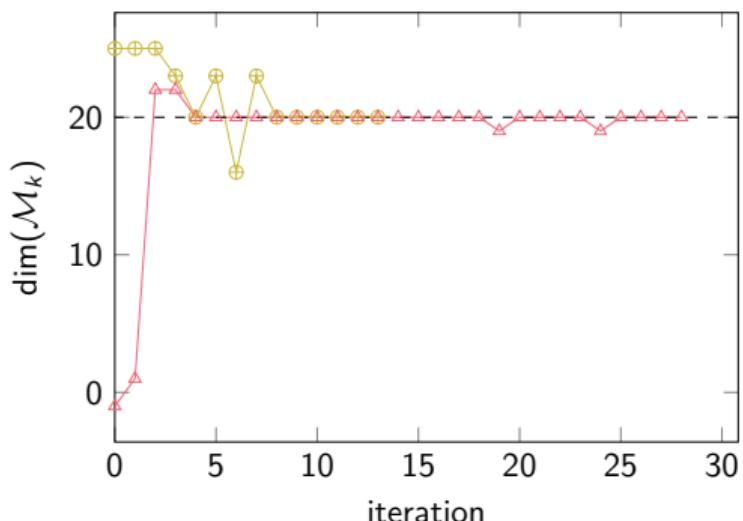
$$\min_{x \in \mathbb{R}^n} \lambda_{\max} \left(A_0 + \sum_{i=1}^n x_i A_i \right)$$

In this instance, $n = 25$, $A_i \in \mathbb{S}_{50}$

$$\mathcal{M}^* = \{x : \lambda_{\max}(c(x)) \text{ has multiplicity 3}\}$$



(our implementation of ◊ Mifflin Sagastizábal '05)



Outline

Introduction

Additive nonsmoothness $f + g$

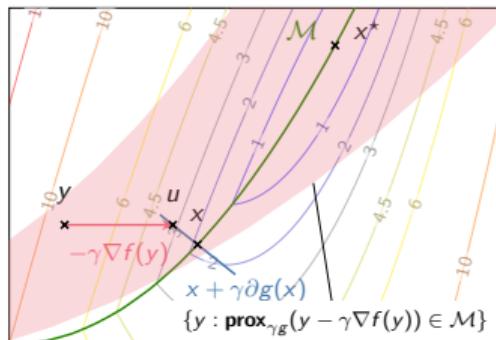
Composite nonsmoothness $g \circ c$

Conclusion

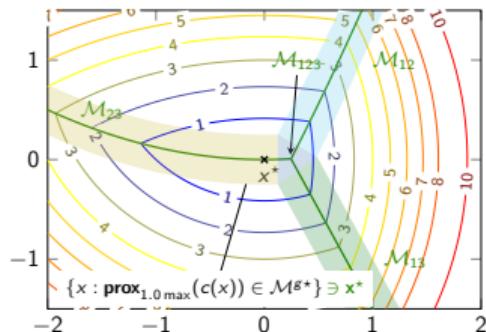
Three years, Three years, in one slide

We have made some steps towards

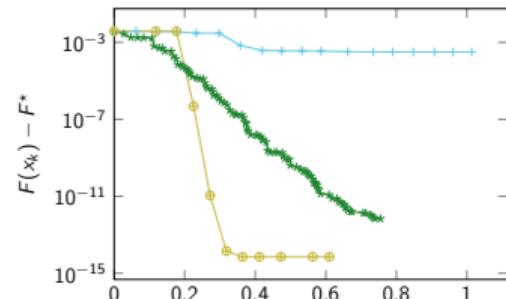
detecting and **exploiting** structure in nonsmooth optimization



Proximal gradient identification



Proximal identification



Quadratic convergence
on nonsmooth functions

Other aspects

Reproducibility

All experiments have online open source implementation in Julia

gbareilles.fr/software

Not discussed but important & interesting

- ▶ active-set solver on spectraplex
- ▶ implementation of \mathcal{VU} -bundle algorithm
- ▶ second derivatives of $\|\cdot\|_*$ and λ_{\max} near nonsmooth points!
- ▶ nonconvex setting

Perspectives

Extensions

- ▶ Globalize local Newton method for $g \circ c$ – work in progress
- ▶ Smarter second-order updates: large-scale settings, efficient even away from minimizers
→ Investigate quasi-Newton (BFGS), cubic regularization of Newton
- ▶ Benchmarking: further numerical comparison ◇ Massias et al '22

Perspectives

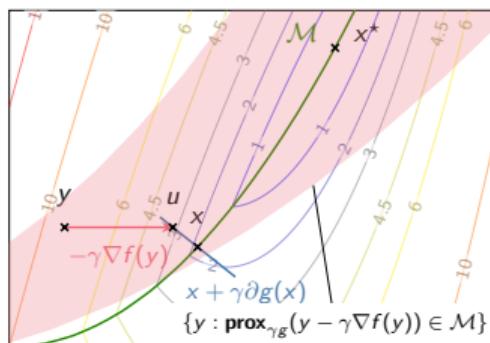
Extensions

- ▶ Globalize local Newton method for $g \circ c$ – work in progress
- ▶ Smarter second-order updates: large-scale settings, efficient even away from minimizers
→ Investigate quasi-Newton (BFGS), cubic regularization of Newton
- ▶ Benchmarking: further numerical comparison ◇ Massias et al '22

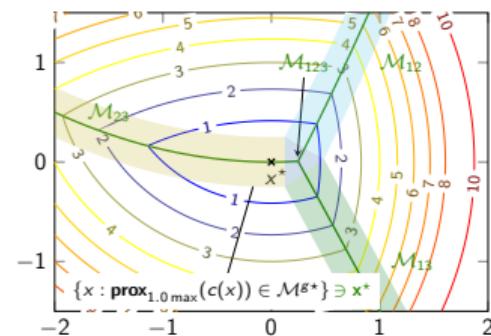
Applications

- ▶ Richer additive structures $f + g + h(L\cdot)$, write as $g \circ c$ with g prox-simple
→ Apply Local Newton Algorithm
- ▶ Smooth minimization on structured sets: $\min_{x \in \mathcal{S}} f(x) \rightarrow \min_x f(x) + \iota_{\mathcal{S}}(x)$
→ Newton acceleration of projected gradient

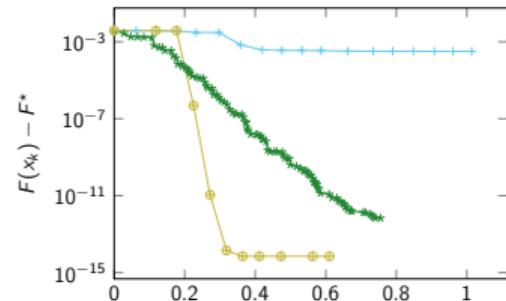
Thank you for your *presence* and *attention!*



Proximal gradient identification



Proximal identification



Quadratic rates
on nonsmooth functions