

Newton methods for nonsmooth composite optimization

Gilles Bareilles

LJK, Univ. Grenoble Alpes

Journées MOA 2022

October 11-14, 2022

Introduction

Detecting structure

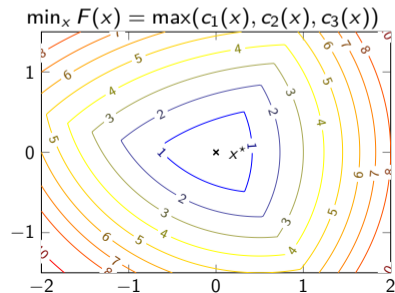
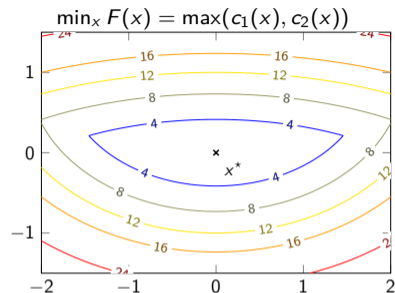
Exploiting structure

Numerics

Conclusion

Composite optimization: $F(x) = g(c(x))$

Includes: max. of \mathcal{C}^2 functions, max. eigenvalue

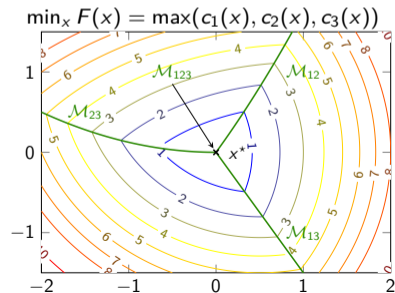
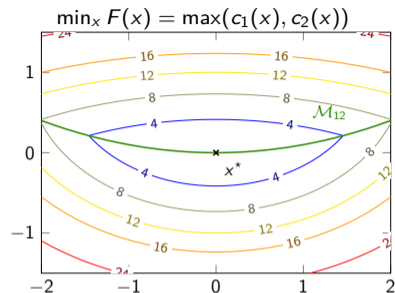


Composite optimization: $F(x) = g(c(x))$

Includes: max. of \mathcal{C}^2 functions, max. eigenvalue

Observations

- nondifferentiability points organize in smooth manifolds
- F is smooth on them



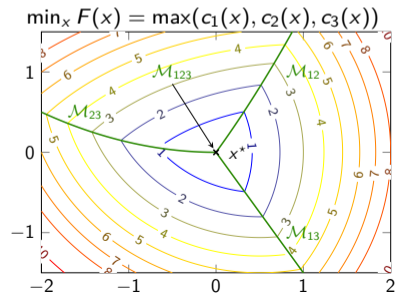
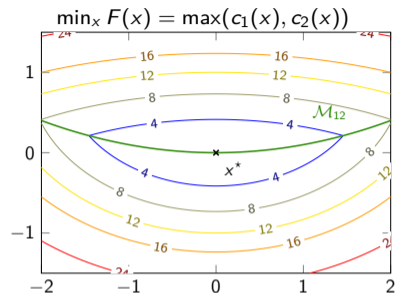
Composite optimization: $F(x) = g(c(x))$

Includes: max. of \mathcal{C}^2 functions, max. eigenvalue

Observations

- nondifferentiability points organize in smooth manifolds
- F is smooth on them

These are **structure manifolds**. ◇ Lewis '02



Composite optimization: $F(x) = g(c(x))$

Includes: max. of \mathcal{C}^2 functions, max. eigenvalue

Observations

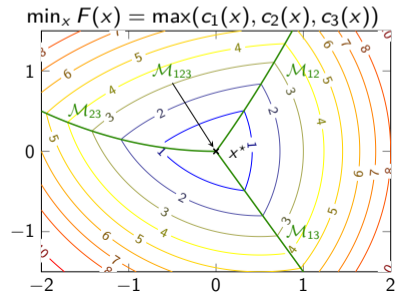
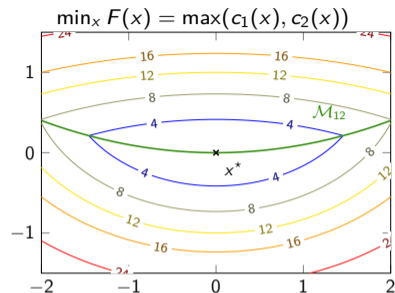
- ▶ nondifferentiability points organize in smooth manifolds
- ▶ F is smooth on them

These are **structure manifolds**. ◇ Lewis '02

Many algorithms for nonsmooth (composite) optimization:

- ▶ prox-linear methods ◇ Lewis Wright, '16,
- ▶ bundle methods ◇ Mifflin Sagastizábal, '05,
- ▶ gradient sampling ◇ Burke Lewis Overton, '05,
- ▶ nonsmooth BFGS ◇ Lewis Overton, '13

Most algorithms are oblivious to structure, **we try to leverage it**.



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

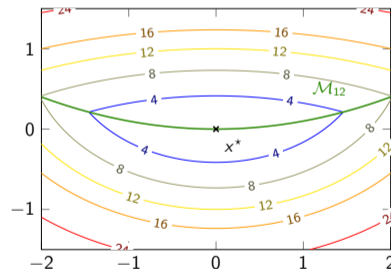
Finding a minimizer of F nonsmooth can be seen as:

- ▶ find the right structure
e.g. which c_i are maximum
- ▶ leverage the right structure to minimize F
e.g. solve smooth problem with smooth constraints

→ We replace (nonsmooth) minimization by smooth constrained minimization.

Challenges:

1. How to detect the optimal structure $\mathcal{M}^* \ni x^*$?
2. How to exploit structure to better minimize F ?



Introduction

Detecting structure

Exploiting structure

Numerics

Conclusion

Prox. for finding structure

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

For *simple functions*, the proximity operator can be computed exactly

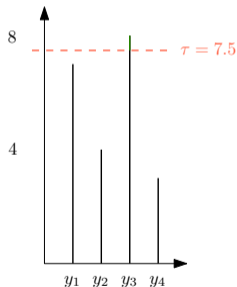
Example (Prox of max)

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

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

Structure manifold:

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



$\gamma = 0.5$

$$\text{prox}_{\gamma \max}(y) = (7, 4, \tau, 3)$$

Structure: \mathcal{M}_I with $I = \{3\}$

Prox. for finding structure

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

For *simple functions*, the proximity operator can be computed exactly

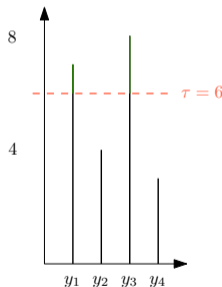
Example (Prox of max)

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

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

Structure manifold:

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



$\gamma = 5$

$$\mathbf{prox}_{\gamma \max}(y) = (\tau, 4, \tau, 3)$$

Structure: \mathcal{M}_I with $I = \{1, 3\}$

Prox. for finding structure

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

For *simple functions*, the proximity operator can be computed exactly

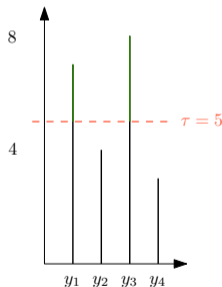
Example (Prox of max)

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

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

Structure manifold:

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



$\gamma = 7$

$$\mathbf{prox}_{\gamma \max}(y) = (\tau, 4, \tau, 3)$$

Structure: \mathcal{M}_I with $I = \{1, 3\}$

Prox. for finding structure

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

For *simple functions*, the proximity operator can be computed exactly

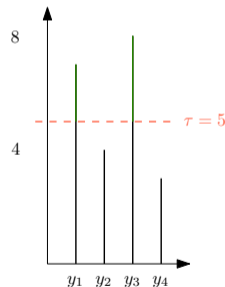
Example (Prox of max)

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

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

Structure manifold:

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



$\gamma = 7$

$$\mathbf{prox}_{\gamma \max}(y) = (\tau, 4, \tau, 3)$$

Structure: \mathcal{M}_I with $I = \{1, 3\}$

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

Identification with explicit prox

Lemma (B., Iutzeler, Malick, '22)

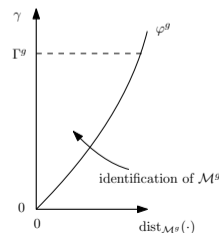
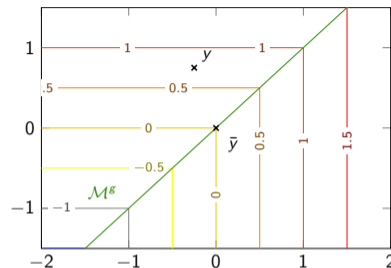
Consider a function g and point \bar{y} with structure \mathcal{M}^g that meet two technical assumptions. For all y near \bar{y} ,

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

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 on the manifold.

Share similarities with [◇ Lewis '02](#), [◇ Lewis Hare '04](#), [◇ Vaiter Peyré Fadili '17](#)



Identification with explicit prox

Lemma (B., Iutzeler, Malick, '22)

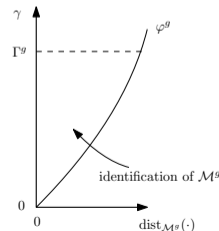
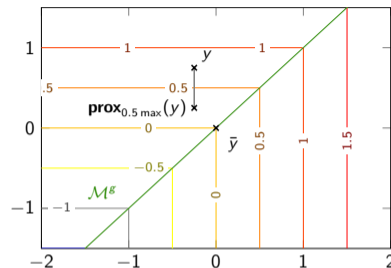
Consider a function g and point \bar{y} with structure \mathcal{M}^g that meet two technical assumptions. For all y near \bar{y} ,

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

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 on the manifold.

Share similarities with [◇ Lewis '02](#), [◇ Lewis Hare '04](#), [◇ Vaiter Peyré Fadili '17](#)



Identification with explicit prox

Lemma (B., Iutzeler, Malick, '22)

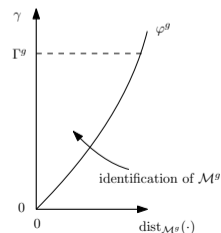
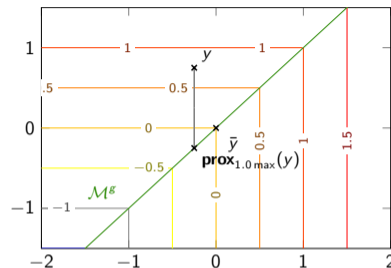
Consider a function g and point \bar{y} with structure \mathcal{M}^g that meet two technical assumptions. For all y near \bar{y} ,

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

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 on the manifold.

Share similarities with [◇ Lewis '02](#), [◇ Lewis Hare '04](#), [◇ Vaïter Peyré Fadili '17](#)



Identification with explicit prox

Lemma (B., Iutzeler, Malick, '22)

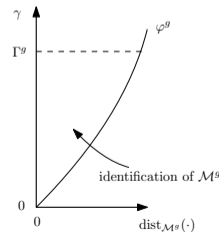
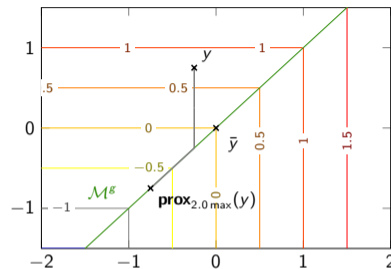
Consider a function g and point \bar{y} with structure \mathcal{M}^g that meet two technical assumptions. For all y near \bar{y} ,

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

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

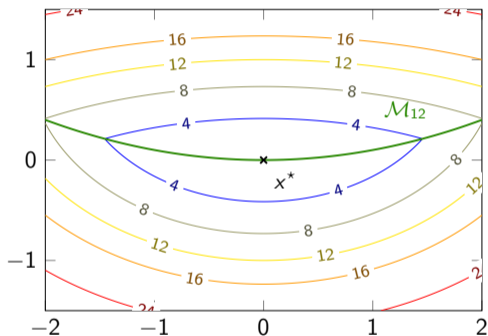
Technical assumptions: normal ascent, control on projection curves on the manifold.

Share similarities with [◇ Lewis '02](#), [◇ Lewis Hare '04](#), [◇ Vaiter Peyré Fadili '17](#)

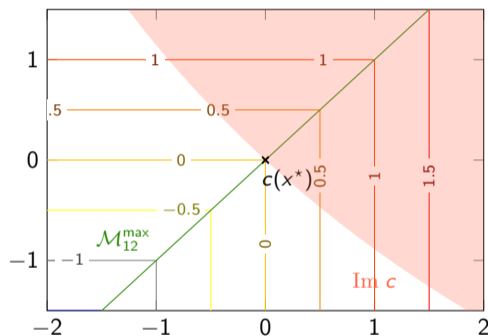


No prox. of F

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



$$F(x) = \max(c_1(x), \dots, c_m(x))$$



$$g(y) = \max(y_1, \dots, y_m)$$

Observation: $\text{prox}_{\gamma g}$ can map points to \mathcal{M}^g .

The structure naturally lies in the intermediate space.

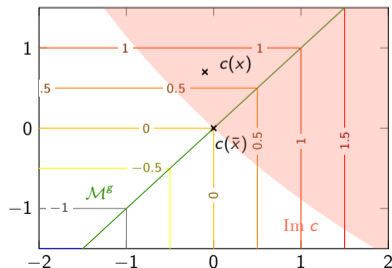
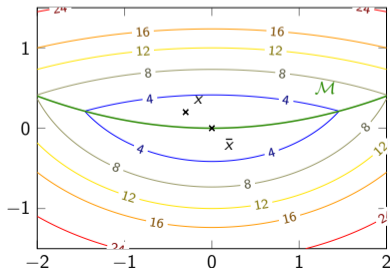
Back to the optimization space

Theorem (B., Iutzeler, Malick, '22)

Consider g , c and a point \bar{x} such that $c(\bar{x})$ has structure manifold \mathcal{M}^g and 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 $\varphi(t) = \frac{c_{\text{map}}}{c_{\text{ri}}} t + \mathcal{O}(t^2)$. Furthermore, $\mathcal{M} = c^{-1}(\mathcal{M}^g)$.



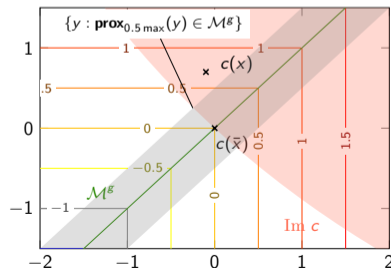
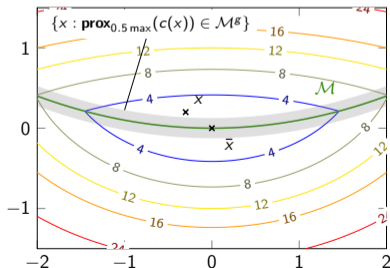
Back to the optimization space

Theorem (B., Iutzeler, Malick, '22)

Consider g , c and a point \bar{x} such that $c(\bar{x})$ has structure manifold \mathcal{M}^g and 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 $\varphi(t) = \frac{c_{\text{map}}}{c_{\text{ri}}} t + \mathcal{O}(t^2)$. Furthermore, $\mathcal{M} = c^{-1}(\mathcal{M}^g)$.



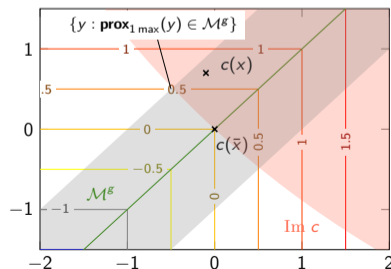
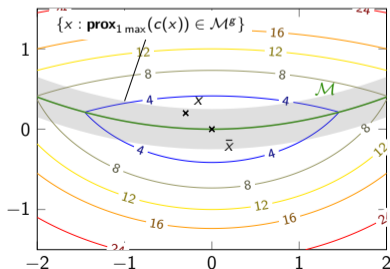
Back to the optimization space

Theorem (B., Iutzeler, Malick, '22)

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

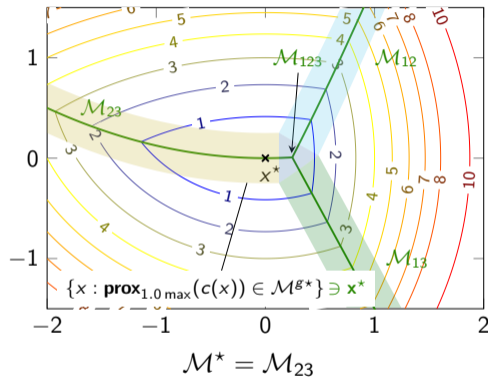
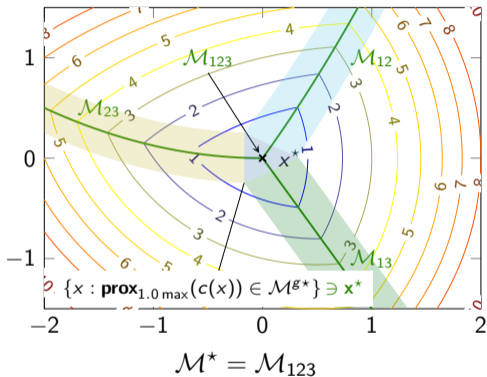
$$\text{prox}_{\gamma g}(c(x)) \in \mathcal{M}^g \quad \text{for all } \gamma \in [\varphi(\text{dist}_{\mathcal{M}}(x)), \Gamma]$$

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



Detection with multiple manifolds

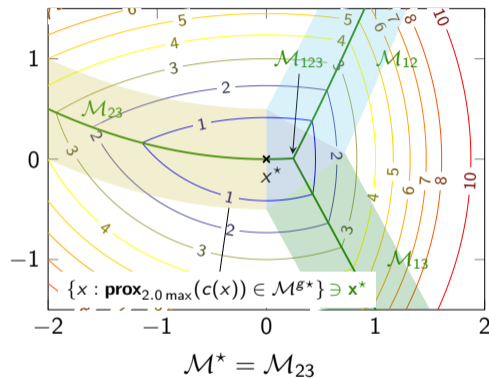
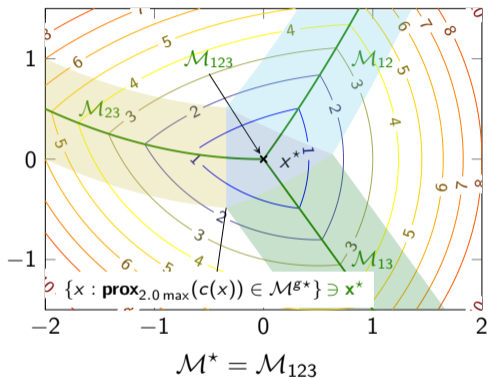
Generally, there are more than one manifolds near x^* .



Importance of γ : too small, detection of \mathcal{M}^* only near x^* ; too large, no detection near x^* .

Detection with multiple manifolds

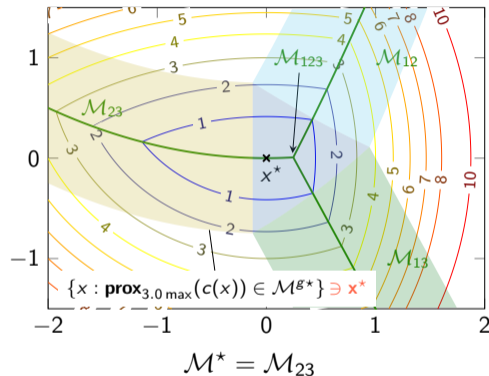
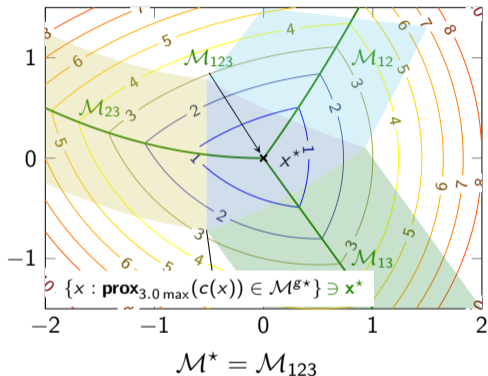
Generally, there are more than one manifolds near x^* .



Importance of γ : too small, detection of \mathcal{M}^* only near x^* ; too large, no detection near x^* .

Detection with multiple manifolds

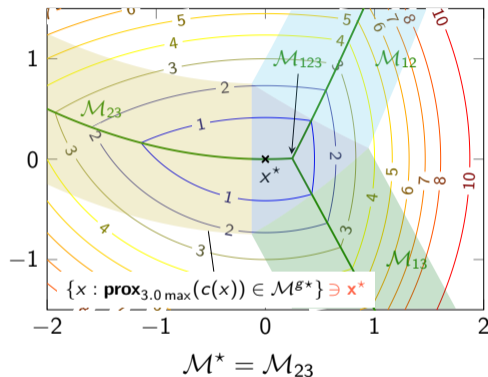
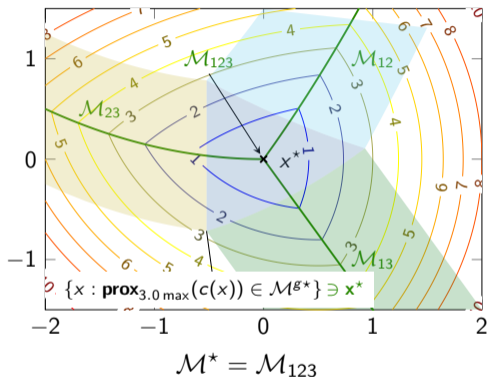
Generally, there are more than one manifolds near x^* .



Importance of γ : too small, detection of \mathcal{M}^* only near x^* ; too large, no detection near x^* .

Detection with multiple manifolds

Generally, there are more than one manifolds near x^* .



Importance of γ : too small, detection of \mathcal{M}^* only near x^* ; too large, no detection near x^* .

Take-away: We detect $\mathcal{M}^* \ni x^*$ with $\text{prox}_{\gamma g} \circ c(\cdot)$ with the right range of steps.

→ How to choose the step in practice?

Introduction

Detecting structure

Exploiting structure

Numerics

Conclusion

Nonsmooth to smooth

- Structure manifolds provide second order models of the nonsmooth F :

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

$$\min_x F(x) \quad \text{and} \quad \mathcal{M} \quad \textbf{turns into} \quad \min_x \tilde{F}(x) \quad \text{s.t. } h(x) = 0.$$

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

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, ...

Newton step and algorithm

Iteration k :

► Compute $\text{prox}_{\gamma_k g}(c(x_k))$ and obtain \mathcal{M}_k .

► With structure candidate \mathcal{M}_k : SQP step on $\min_x \tilde{F}_k(x)$ s.t. $h_k(x) = 0$.

$$d_k^{\text{SQP}} = \arg \min_{d \in \mathbb{R}^n} \quad \langle \nabla \tilde{F}_k(x_k), d \rangle + \frac{1}{2} \langle \nabla_{xx}^2 L_k(x_k, \lambda_k(x_k)) d, d \rangle$$

$$\text{s.t.} \quad h_k(x_k) + D h_k(x_k) d = 0$$

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

► $\gamma_{k+1} = \frac{\gamma_k}{2}$

Similar works with *heuristic* structure detection: ♦ Womersley Fletcher '86 for max, ♦ Noll Apkarian, '05 for λ_{\max} .

Local exact structure identification and quadratic convergence

Theorem (B., Iutzeler, Malick, '22)

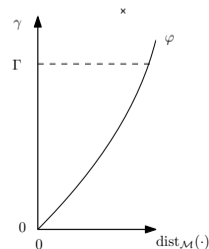
Consider a function $F = g \circ c$ and x^* a strong minimizer with structure manifold \mathcal{M}^* that meets the technical assumptions.

If x_0 and $F(x_0)$ are close enough to x^* and $F(x^*)$, γ_0 is large enough and no Maratos effect happens, then there exists $C > 0$ such that:

$$\mathcal{M}_k = \mathcal{M}^* \quad \text{and} \quad \|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2 \quad \text{for all } k \text{ large enough.}$$

Proof idea

- ▶ if $\mathcal{M}_k = \mathcal{M}^*$, the SQP step brings quadratic improvement
- ▶ since γ_k decreases, at some point $\gamma_k \in [\varphi(\text{dist}_{\mathcal{M}}(x_k)), \Gamma]$
- ▶ to stay in that region, decrease γ not too fast



Local exact structure identification and quadratic convergence

Theorem (B., Iutzeler, Malick, '22)

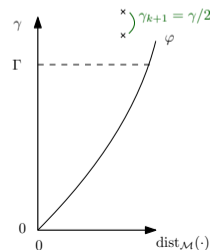
Consider a function $F = g \circ c$ and x^* a strong minimizer with structure manifold \mathcal{M}^* that meets the technical assumptions.

If x_0 and $F(x_0)$ are close enough to x^* and $F(x^*)$, γ_0 is large enough and no Maratos effect happens, then there exists $C > 0$ such that:

$$\mathcal{M}_k = \mathcal{M}^* \quad \text{and} \quad \|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2 \quad \text{for all } k \text{ large enough.}$$

Proof idea

- ▶ if $\mathcal{M}_k = \mathcal{M}^*$, the SQP step brings quadratic improvement
- ▶ since γ_k decreases, at some point $\gamma_k \in [\varphi(\text{dist}_{\mathcal{M}}(x_k)), \Gamma]$
- ▶ to stay in that region, decrease γ not too fast



Local exact structure identification and quadratic convergence

Theorem (B., Iutzeler, Malick, '22)

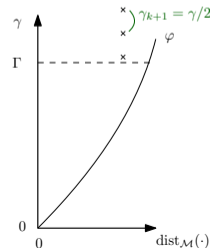
Consider a function $F = g \circ c$ and x^* a strong minimizer with structure manifold \mathcal{M}^* that meets the technical assumptions.

If x_0 and $F(x_0)$ are close enough to x^* and $F(x^*)$, γ_0 is large enough and no Maratos effect happens, then there exists $C > 0$ such that:

$$\mathcal{M}_k = \mathcal{M}^* \quad \text{and} \quad \|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2 \quad \text{for all } k \text{ large enough.}$$

Proof idea

- ▶ if $\mathcal{M}_k = \mathcal{M}^*$, the SQP step brings quadratic improvement
- ▶ since γ_k decreases, at some point $\gamma_k \in [\varphi(\text{dist}_{\mathcal{M}}(x_k)), \Gamma]$
- ▶ to stay in that region, decrease γ not too fast



Local exact structure identification and quadratic convergence

Theorem (B., Iutzeler, Malick, '22)

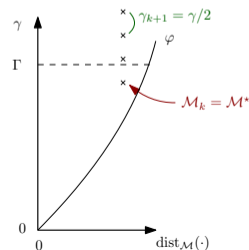
Consider a function $F = g \circ c$ and x^* a strong minimizer with structure manifold \mathcal{M}^* that meets the technical assumptions.

If x_0 and $F(x_0)$ are close enough to x^* and $F(x^*)$, γ_0 is large enough and no Maratos effect happens, then there exists $C > 0$ such that:

$$\mathcal{M}_k = \mathcal{M}^* \quad \text{and} \quad \|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2 \quad \text{for all } k \text{ large enough.}$$

Proof idea

- ▶ if $\mathcal{M}_k = \mathcal{M}^*$, the SQP step brings quadratic improvement
- ▶ since γ_k decreases, at some point $\gamma_k \in [\varphi(\text{dist}_{\mathcal{M}}(x_k)), \Gamma]$
- ▶ to stay in that region, decrease γ not too fast



Local exact structure identification and quadratic convergence

Theorem (B., Iutzeler, Malick, '22)

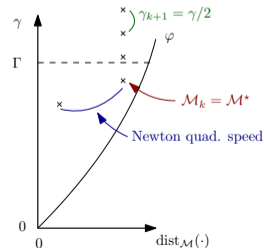
Consider a function $F = g \circ c$ and x^* a strong minimizer with structure manifold \mathcal{M}^* that meets the technical assumptions.

If x_0 and $F(x_0)$ are close enough to x^* and $F(x^*)$, γ_0 is large enough and no Maratos effect happens, then there exists $C > 0$ such that:

$$\mathcal{M}_k = \mathcal{M}^* \quad \text{and} \quad \|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2 \quad \text{for all } k \text{ large enough.}$$

Proof idea

- ▶ if $\mathcal{M}_k = \mathcal{M}^*$, the SQP step brings quadratic improvement
- ▶ since γ_k decreases, at some point $\gamma_k \in [\varphi(\text{dist}_{\mathcal{M}}(x_k)), \Gamma]$
- ▶ to stay in that region, decrease γ not too fast



Local exact structure identification and quadratic convergence

Theorem (B., Iutzeler, Malick, '22)

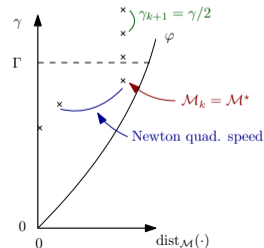
Consider a function $F = g \circ c$ and x^* a strong minimizer with structure manifold \mathcal{M}^* that meets the technical assumptions.

If x_0 and $F(x_0)$ are close enough to x^* and $F(x^*)$, γ_0 is large enough and no Maratos effect happens, then there exists $C > 0$ such that:

$$\mathcal{M}_k = \mathcal{M}^* \quad \text{and} \quad \|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2 \quad \text{for all } k \text{ large enough.}$$

Proof idea

- ▶ if $\mathcal{M}_k = \mathcal{M}^*$, the SQP step brings quadratic improvement
- ▶ since γ_k decreases, at some point $\gamma_k \in [\varphi(\text{dist}_{\mathcal{M}}(x_k)), \Gamma]$
- ▶ to stay in that region, decrease γ not too fast



Local exact structure identification and quadratic convergence

Theorem (B., Iutzeler, Malick, '22)

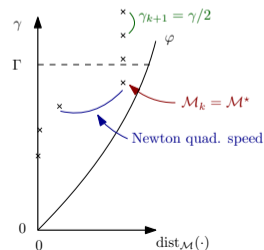
Consider a function $F = g \circ c$ and x^* a strong minimizer with structure manifold \mathcal{M}^* that meets the technical assumptions.

If x_0 and $F(x_0)$ are close enough to x^* and $F(x^*)$, γ_0 is large enough and no Maratos effect happens, then there exists $C > 0$ such that:

$$\mathcal{M}_k = \mathcal{M}^* \quad \text{and} \quad \|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2 \quad \text{for all } k \text{ large enough.}$$

Proof idea

- ▶ if $\mathcal{M}_k = \mathcal{M}^*$, the SQP step brings quadratic improvement
- ▶ since γ_k decreases, at some point $\gamma_k \in [\varphi(\text{dist}_{\mathcal{M}}(x_k)), \Gamma]$
- ▶ to stay in that region, decrease γ not too fast



Introduction

Detecting structure

Exploiting structure

Numerics

Conclusion

Quadratic convergence

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

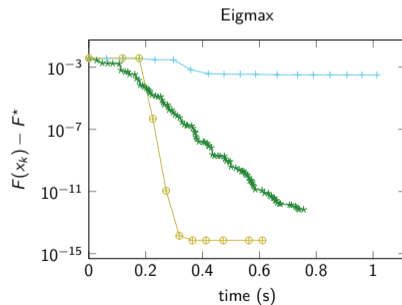
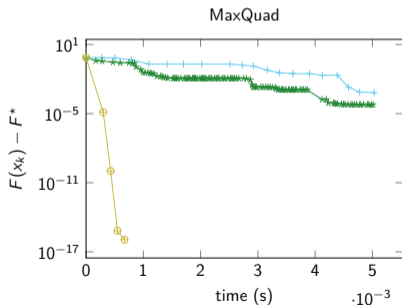
$$\mathcal{M}^* = \{x : c_2(x) = \dots = c_5(x)\}$$

$$\min_{x \in \mathbb{R}^{25}} \lambda_{\max} \left(A_0 + \sum_{i=1}^n x_i A_i \right)$$

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

Historical maxquad problem ◇ HULL '93

Matrices are symmetric, 50×50



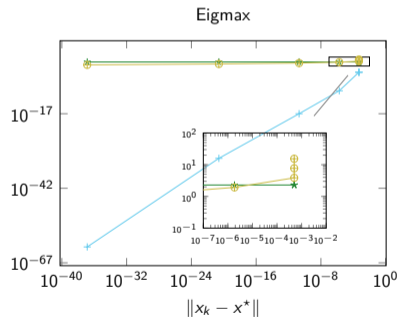
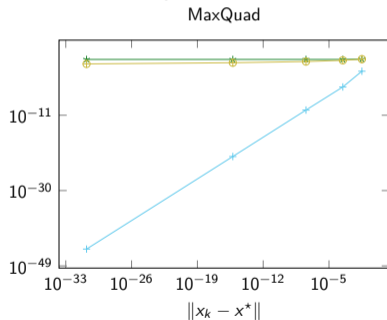
—+— Gradient Sampling —*— nsBFGS —○— LocalNewton

Proximal identification

Corollary: There exists $L > 0$, $\epsilon > 0$ such that

$$\|x - x^*\| \leq \epsilon \text{ and } L\|x - x^*\| \leq \gamma \leq \Gamma \implies \text{prox}_{\gamma g}(c(x)) \in \mathcal{M}^{g^*}.$$

This checks out in practice:



—+— $\inf\{\gamma : \text{prox}_{\gamma g}(c(x_k)) \in \mathcal{M}^{g^*}\}$ —*— $\sup\{\gamma : \text{prox}_{\gamma g}(c(x_k)) \in \mathcal{M}^{g^*}\}$ —⊕— γ_k

Introduction

Detecting structure

Exploiting structure

Numerics

Conclusion

Conclusion

Take-away messages

- ▶ Proximal methods identify smooth structure in nonsmooth composite problems
- ▶ We show local **exact** identification and quadratic rate for $g \circ c$, where g is prox-simple, no convexity required

B. & Iutzeler & Malick: Harnessing structure in composite nonsmooth minimization

<https://arxiv.org/abs/2206.15053>

Work in progress and perspectives

- ▶ Drop the locality: i) need more information to identify, ii) globalize constrained Newton

Thank you!

Technical assumptions

Normal ascent: g increases at \bar{y} on normal directions:

$$0 \in \text{ri } \mathbf{proj}_{N_{\bar{y}}\mathcal{M}^g} \partial g(\bar{y})$$

Manifold curves: A function g with structure \mathcal{M}^g at \bar{y} satisfies the *curve property* if there exists a neighborhood $\mathcal{N}_{\bar{y}}$ of \bar{y} and $T > 0$ such that, for any smooth application $e : \mathcal{N}_{\bar{y}} \times [0, T] \rightarrow \mathcal{M}^g$ verifying $e(y, 0) = \mathbf{proj}_{\mathcal{M}^g}(y)$ and $\frac{d}{dt}e(y, 0) = -\text{grad } g(\mathbf{proj}_{\mathcal{M}^g}(y))$, there holds

$$\|\mathbf{proj}_{N_{e(y,t)}\mathcal{M}^g}(e(y, t) - y)\| \leq \text{dist}_{\mathcal{M}^g}(y) + \tilde{L} t^2 \quad \text{for all } y \in \mathcal{N}_{\bar{y}}, t \in [0, T],$$

where $\text{grad } g(p) \in T_p\mathcal{M}^g$ denotes the Riemannian gradient of g , obtained as $\mathbf{proj}_{T_p\mathcal{M}^g}(\partial g(p))$.

No Maratos: near a minimizer x^* , a step d that makes $x + d$ quadratically closer to x^* than x implies descent $F(x + d) \leq F(x)$.

Transversality: the mapping $c : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is transversal to manifold $\mathcal{M} \subset \mathbb{R}^m$ at $c(x)$ if:

$$\ker(\text{Jac}_c(x)^\top) \cap N_{c(x)}\mathcal{M}^g = \{0\}$$

\Rightarrow if $\text{Jac}_h(c(x))$ is full rank, then $\text{Jac}_{h \circ c}(x)$ is also full-rank.

Maximum structure and initial stepsize

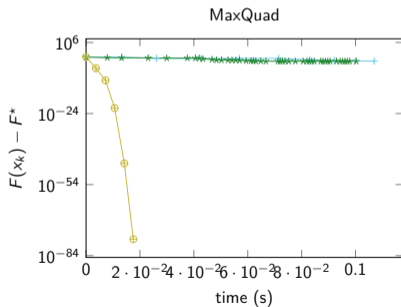
In the generated instance, the multiplicity of the maximum eigenvalue at optimum is $r = 3$. The maximum structure of a point, useful in setting γ_0 , is \mathcal{M}_r , with $r = 6$, and not the matrix size $m = 50$. Indeed, the codimension of \mathcal{M}_r , that is the dimension of its normal spaces, should be lower than that of \mathbb{R}^n : $r(r+1)/2 - 1 \leq 25$, that is $r \leq 6$ (see the discussion in [?, pp. 555-556, Eq. 2.5]).

Quadratic convergence, BigFloat precision

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

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

Historical maxquad problem



$$\min_{x \in \mathbb{R}^{25}} \lambda_{\max} \left(A_0 + \sum_{i=1}^n x_i A_i \right)$$

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

Matrices are symmetric, 50×50 .

