# Smoothness in nonsmooth optimization

(Newtonian ideas for partly smooth equations)

**Adrian Lewis** 

Joint work with: D. Drusvyatskiy, X.Y. Han, A. loffe, J. Liang, M.L. Overton, T. Tian, C. Wylie

September 2020

**ORIE** Cornell

One World Optimization Seminar

**Question 2:** Can Newton methods use this smoothness?

**Question 2:** Can Newton methods use this smoothness?

**Example:**  $\min_Q f$  becomes  $-\nabla f(x) \in N_Q(x)$ . Projected gradient methods  $x \leftarrow \operatorname{Proj}_Q (x - \gamma \nabla f(x))$ identify smoothness in Q. Newtonian acceleration?

**Question 2:** Can Newton methods use this smoothness?

**Example:**  $\min_Q f$  becomes  $-\nabla f(x) \in N_Q(x)$ . Projected gradient methods  $x \leftarrow \operatorname{Proj}_Q (x - \gamma \nabla f(x))$ identify smoothness in Q. Newtonian acceleration? **Question 3:** Superlinear convergence for black box nonsmooth optimization?

**Question 2:** Can Newton methods use this smoothness?

**Example:**  $\min_Q f$  becomes  $-\nabla f(x) \in N_Q(x)$ . Projected gradient methods  $x \leftarrow \operatorname{Proj}_Q (x - \gamma \nabla f(x))$ identify smoothness in Q. Newtonian acceleration? **Question 3:** Superlinear convergence for black box nonsmooth optimization?

**Example:** Eigenvalue optimization.



#### Inherent structure: an example

The numerical radius of an *n*-by-*n* complex matrix *A*,

$$\rho(A) = \max_{\|u\|=1} |u^*Au|,$$

satisfies the "power inequality" (Berger '65): for k = 1, 2, ...

$$\frac{1}{2}\|\boldsymbol{A}^{k}\|_{2} \leq \rho(\boldsymbol{A}^{k}) \leq (\rho(\boldsymbol{A}))^{k},$$

and so controls transient stability in dynamics  $x \leftarrow Ax$ .

The numerical radius of an *n*-by-*n* complex matrix *A*,

$$\rho(A) = \max_{\|u\|=1} |u^*Au|,$$

satisfies the "power inequality" (Berger '65): for k = 1, 2, ...

$$\frac{1}{2}\|\boldsymbol{A}^{k}\|_{2} \leq \rho(\boldsymbol{A}^{k}) \leq (\rho(\boldsymbol{A}))^{k},$$

and so controls transient stability in dynamics  $x \leftarrow Ax$ .

Optimizing  $\rho$  often results in unusual matrices...

The numerical radius of an *n*-by-*n* complex matrix *A*,

$$\rho(A) = \max_{\|u\|=1} |u^*Au|,$$

satisfies the "power inequality" (Berger '65): for k = 1, 2, ...

$$\frac{1}{2}\|A^k\|_2 \leq \rho(A^k) \leq (\rho(A))^k,$$

and so controls transient stability in dynamics  $x \leftarrow Ax$ .

Optimizing  $\rho$  often results in unusual matrices...

Example: For random matrices Y, proximal matrices A minimizing

$$\rho(A) + \lambda \|A - Y\|^2$$

often have disc fields of values  $\{u^*Au : ||u|| = 1\}$ .

#### Random proximal points (via cvx) are often disk matrices

$$1 - rac{ ext{inner radius}}{
ho(A)}$$
(via chebfun)

algebraic deviation from disk

distance to

singularity

distance to null Jordan block

### Random proximal points (via cvx) are often disk matrices

 $1 - rac{\text{inner radius}}{\rho(A)}$ (via chebfun)

algebraic deviation from disk

distance to singularity

distance to null Jordan block



### Random proximal points (via cvx) are often disk matrices

 $1 - \frac{\text{inner radius}}{\rho(A)}$ (via chebfun)

algebraic deviation from disk

distance to singularity

distance to null Jordan block



**Why...?** (Disk matrices comprise a small set, of codimension 2n).  $_{3/24}$ 

• {smooth  $g_i(x) \le 0$ } relative to active set  $\{x : g_j(x) = 0\}$ 

L '02

- {smooth  $g_i(x) \le 0$ } relative to active set  $\{x : g_j(x) = 0\}$
- PSD matrices  $\mathbf{S}_{+}^{n}$  relative to  $\{X \in \mathbf{S}_{+}^{n} : \operatorname{rank}(X) = k\}$

- {smooth  $g_i(x) \le 0$ } relative to active set  $\{x : g_j(x) = 0\}$
- PSD matrices  $\mathbf{S}_{+}^{n}$  relative to  $\{X \in \mathbf{S}_{+}^{n} : \operatorname{rank}(X) = k\}$
- smooth  $f + || \cdot ||_1$  relative to fixed sparsity pattern

- {smooth  $g_i(x) \le 0$ } relative to active set  $\{x : g_j(x) = 0\}$
- PSD matrices  $\mathbf{S}_{+}^{n}$  relative to  $\{X \in \mathbf{S}_{+}^{n} : \operatorname{rank}(X) = k\}$
- smooth  $f + || \cdot ||_1$  relative to fixed sparsity pattern
- numerical radius  $\rho$  relative to disk matrices (L-Overton '20, Han-L '20)

- {smooth  $g_i(x) \le 0$ } relative to active set  $\{x : g_j(x) = 0\}$
- PSD matrices  $\mathbf{S}_{+}^{n}$  relative to  $\{X \in \mathbf{S}_{+}^{n} : \operatorname{rank}(X) = k\}$
- smooth  $f + || \cdot ||_1$  relative to fixed sparsity pattern
- numerical radius ρ relative to disk matrices (L-Overton '20, Han-L '20)

**History:** generalized "active constraints" in nonlinear programming (Burke-Moré '88), "identifiable surfaces" (Wright '93), "VU decomposition" (Mifflin-Sagastizábal '00)...

Diverse first-order methods identify the manifold (L-Hare '04...), which drives the local convergence.

Diverse first-order methods identify the manifold (L-Hare '04...), which drives the local convergence.

**Example:** *l*<sub>1</sub> regularization for sparsity. Proximal gradient iterates settle on a sparsity pattern (Hale-Yin-Zhang '08).

Diverse first-order methods identify the manifold (L-Hare '04...), which drives the local convergence.

**Example:** *l*<sub>1</sub> regularization for sparsity. Proximal gradient iterates settle on a sparsity pattern (Hale-Yin-Zhang '08).



$$\min_{x} f(x) = h(x) + \lambda \|x\|_*$$

Diverse first-order methods identify the manifold (L-Hare '04...), which drives the local convergence.

**Example:** *l*<sub>1</sub> regularization for sparsity. Proximal gradient iterates settle on a sparsity pattern (Hale-Yin-Zhang '08).



$$\min_{x} f(x) = h(x) + \lambda \|x\|_{*}$$

**Example:** Nuclear norm regularization for low-rank optimization. Proximal gradient (singular value thresholding) iterates settle on a fixed-rank manifold, then converge linearly to the solution.

(Liang-Fadili-Peyré '18)

Shift focus from optimization to optimality conditions:

x minimizes  $f \Rightarrow 0 \in \partial f(x)$ .

Shift focus from optimization to optimality conditions:

x minimizes  $f \Rightarrow 0 \in \partial f(x)$ .

Generalize:  $0 \in \Phi(u)$  for set-valued operator  $\Phi$  on  $\mathbb{R}^n$ .

Shift focus from optimization to optimality conditions:

x minimizes 
$$f \Rightarrow 0 \in \partial f(x)$$
.

Generalize:  $0 \in \Phi(u)$  for set-valued operator  $\Phi$  on  $\mathbb{R}^n$ .

• Variational inequalities

Find  $x \in Q$  so  $F(x)^T(z-x) \ge 0$  for all  $z \in Q$ : equivalently,

 $0\in F(x)+N_Q(x).$ 

Shift focus from optimization to optimality conditions:

x minimizes 
$$f \Rightarrow 0 \in \partial f(x)$$
.

Generalize:  $0 \in \Phi(u)$  for set-valued operator  $\Phi$  on  $\mathbb{R}^n$ .

• Variational inequalities Find  $x \in Q$  so  $F(x)^T(z - x) \ge 0$  for all  $z \in Q$ : equivalently,  $0 \in F(x) + N_Q(x)$ .

#### • Composite optimization

 $\min_{x} h(c(x))$  for convex h on  $\mathbb{R}^{m}$  and smooth c into  $\mathbb{R}^{m}$ . Stationarity:

$$0 \in \left( \nabla c(x)^T y \\ -y \right) + \left( \begin{matrix} 0 \\ \partial h(c(x)) \end{matrix} \right)$$

## Partly smooth generalized equations (L-Liang '18)

A set-valued operator  $\Phi$  is partly smooth at a solution  $\bar{u}$  for given data  $\bar{v}$ 

# Partly smooth generalized equations (L-Liang '18)

A set-valued operator  $\Phi$  is partly smooth at a solution  $\bar{u}$  for given data  $\bar{v}$  if

- gph  $\Phi = \{(u, v) : v \in \Phi(u)\}$ is a manifold around  $(\overline{u}, \overline{v})$ .
- proj : gph Φ → R<sup>n</sup> : (u, v) → u
   is constant rank around (ū, v)...
   (i.e. the projection of the graph's tangent space has locally constant dimension).



## Partly smooth generalized equations (L-Liang '18)

A set-valued operator  $\Phi$  is partly smooth at a solution  $\bar{u}$  for given data  $\bar{v}$  if

- gph  $\Phi = \{(u, v) : v \in \Phi(u)\}$ is a manifold around  $(\overline{u}, \overline{v})$ .
- proj : gph Φ → R<sup>n</sup> : (u, v) → u
   is constant rank around (ū, v)...
   (i.e. the projection of the graph's tangent space has locally constant dimension).



Asymptotic solvers then identify the active manifold

$$\mathcal{M} \;=\; {
m proj}ig({
m gph}\,\Phi\,\,{
m around}\,\,(ar{u},ar{v})ig),$$

since  $v_k \in \Phi(u_k)$  with  $(u_k, v_k) \to (\bar{u}, \bar{v})$  implies  $u_k \in \mathcal{M}$  eventually.

#### Basic example: partly smooth sets

For closed convex (or "prox-regular")  $S \subset \mathbb{R}^n$ , suppose  $\bar{x}$  solves  $\min_{x \in S} \langle \bar{y}, x \rangle \quad \text{and hence} \quad \bar{y} \in N_S(x).$ 

8/24

#### Basic example: partly smooth sets

For closed convex (or "prox-regular")  $S \subset \mathbf{R}^n$ , suppose  $\bar{x}$  solves  $\min_{x \in S} \langle \bar{y}, x \rangle$  and hence  $\bar{y} \in N_S(x)$ .

If S contains a ridge manifold  $\mathcal{M}$ (the normal cone  $N_S(x)$  depends on  $x \in \mathcal{M}$  continuously and spans  $N_{\mathcal{M}}(x)$ ), and nondegeneracy holds  $(\bar{y} \in ri(N_S(\bar{x})))$ , then the operator  $N_S$  is partly smooth at  $\bar{x}$  for  $\bar{y}$ , with active manifold  $\mathcal{M}$ .



#### Basic example: partly smooth sets

For closed convex (or "prox-regular")  $S \subset \mathbf{R}^n$ , suppose  $\bar{x}$  solves  $\min_{x \in S} \langle \bar{y}, x \rangle$  and hence  $\bar{y} \in N_S(x)$ .

If S contains a ridge manifold  $\mathcal{M}$ (the normal cone  $N_S(x)$  depends on  $x \in \mathcal{M}$  continuously and spans  $N_{\mathcal{M}}(x)$ ), and nondegeneracy holds  $(\bar{y} \in ri(N_S(\bar{x})))$ , then the operator  $N_S$  is partly smooth at  $\bar{x}$  for  $\bar{y}$ , with active manifold  $\mathcal{M}$ .



So if S is convex and  $\bar{x}$  is unique, projected gradient iterations  $x \leftarrow \operatorname{proj}_Q(x - \alpha \bar{y})$  converge to  $\bar{x}$  (if  $\alpha$  small) and identify  $\mathcal{M}$ .
#### **Example:** max functions of degree *k*

### **Example:** max functions of degree k

Given a decomposition

$$f(x) = \max_{i=1,\ldots,k} f_i(x),$$

using smooth components  $f_i$ ,

Given a decomposition

$$f(x) = \max_{i=1,\ldots,k} f_i(x),$$

using smooth components  $f_i$ , call  $\bar{x}$  a strictly active critical point when the values  $f_i(\bar{x})$  are all equal, and the system

$$\sum_{i} \lambda_{i} = 1, \qquad \sum_{i} \lambda_{i} \nabla f_{i}(\bar{x}) = 0$$

has a unique solution, which furthermore has each  $\lambda_i > 0$ .

Given a decomposition

$$f(x) = \max_{i=1,\ldots,k} f_i(x),$$

using smooth components  $f_i$ , call  $\bar{x}$  a strictly active critical point when the values  $f_i(\bar{x})$  are all equal, and the system

$$\sum_{i} \lambda_{i} = 1, \qquad \sum_{i} \lambda_{i} \nabla f_{i}(\bar{x}) = 0$$

has a unique solution, which furthermore has each  $\lambda_i > 0$ . Then

$$x \mapsto \partial f(x) = \operatorname{conv}\{\nabla f_i(x) : f_i(x) = f(x)\}$$

is partly smooth at  $\bar{x}$  for 0 relative to the active manifold  $\mathcal{M}$  of points where each  $f_i$  has equal value.

Sard's Theorem: almost no values of smooth operators are critical.

Sard's Theorem: almost no values of smooth operators are critical. What about set-valued operators and generalized equations on  $\mathbb{R}^n$ ?

Sard's Theorem: almost no values of smooth operators are critical. What about set-valued operators and generalized equations on  $\mathbb{R}^n$ ?

Consider a semi-algebraic operator  $\Phi$  with *n*-dimensional graph:

$$gph \Phi = \bigcup_{i=1}^{q} \bigcap_{j=1}^{r} \left\{ (x, y) \in \mathbf{R}^{2n} : p_{ij}(x, y) \left\{ \begin{array}{c} \leq \\ \text{or} \\ < \end{array} \right\} = 0 \right\}$$

for polynomials  $p_{ij}$ .

Sard's Theorem: almost no values of smooth operators are critical. What about set-valued operators and generalized equations on  $\mathbb{R}^n$ ?

Consider a semi-algebraic operator  $\Phi$  with *n*-dimensional graph:

$$gph \Phi = \bigcup_{i=1}^{q} \bigcap_{j=1}^{r} \left\{ (x, y) \in \mathbf{R}^{2n} : p_{ij}(x, y) \left\{ \begin{array}{c} \leq \\ \text{or} \\ < \end{array} \right\} 0 \right\}$$

for polynomials p<sub>ij</sub>. Eg: Subdifferentials, monotone operators...

Sard's Theorem: almost no values of smooth operators are critical. What about set-valued operators and generalized equations on  $\mathbb{R}^{n}$ ?

Consider a semi-algebraic operator  $\Phi$  with *n*-dimensional graph:

$$gph \Phi = \bigcup_{i=1}^{q} \bigcap_{j=1}^{r} \left\{ (x, y) \in \mathbf{R}^{2n} : p_{ij}(x, y) \left\{ \begin{array}{c} \leq \\ \text{or} \\ < \end{array} \right\} 0 \right\}$$

for polynomials p<sub>ij</sub>. Eg: Subdifferentials, monotone operators...

**Theorem** Around generic data y, there are smooth maps  $G_i$  so  $\Phi^{-1} = \{G_1, \dots, G_k\}$  (possibly empty).

Sard's Theorem: almost no values of smooth operators are critical. What about set-valued operators and generalized equations on  $\mathbb{R}^n$ ?

Consider a semi-algebraic operator  $\Phi$  with *n*-dimensional graph:

$$gph \Phi = \bigcup_{i=1}^{q} \bigcap_{j=1}^{r} \left\{ (x, y) \in \mathbf{R}^{2n} : p_{ij}(x, y) \left\{ \begin{array}{c} \leq \\ \text{or} \\ < \end{array} \right\} 0 \right\}$$

for polynomials p<sub>ij</sub>. Eg: Subdifferentials, monotone operators...

**Theorem** Around generic data y, there are smooth maps  $G_i$  so  $\Phi^{-1} = \{G_1, \dots, G_k\}$  (possibly empty).

 $\Phi$  is partly smooth for y at each solution  $x_i = G_i(y)$ ,

Sard's Theorem: almost no values of smooth operators are critical. What about set-valued operators and generalized equations on  $\mathbb{R}^{n}$ ?

Consider a semi-algebraic operator  $\Phi$  with *n*-dimensional graph:

$$gph \Phi = \bigcup_{i=1}^{q} \bigcap_{j=1}^{r} \left\{ (x, y) \in \mathbf{R}^{2n} : p_{ij}(x, y) \left\{ \begin{array}{c} \leq \\ \text{or} \\ < \end{array} \right\} = 0 \right\}$$

for polynomials  $p_{ij}$ . Eg: Subdifferentials, monotone operators...

**Theorem** Around generic data y, there are smooth maps  $G_i$  so  $\Phi^{-1} = \{G_1, \dots, G_k\}$  (possibly empty).

 $\Phi$  is partly smooth for y at each solution  $x_i = G_i(y)$ , and

gph  $\Phi$  intersects ( $\mathbb{R}^n \times \{y\}$ ) transversally at  $(x_i, y)$ . (loffe '07, Bolte...'11, Drusvyatskiy...'16, Lee...'19, L-Tian)

Recast as set intersection: find a point z = (u, 0) where

 $X = \operatorname{gph} \Phi$  intersects  $Y = \mathbb{R}^n \times \{0\}.$ 

# (L-Wylie '20)

Recast as set intersection: find a point z = (u, 0) where

 $X = \operatorname{gph} \Phi$  intersects  $Y = \mathbf{R}^n \times \{0\}.$ 

Assume transversality:  $N_X(z) \cap N_Y(z) = \{0\}.$ 

Recast as set intersection: find a point z = (u, 0) where

 $X = \operatorname{gph} \Phi$  intersects  $Y = \mathbb{R}^n \times \{0\}.$ 

Assume transversality:  $N_X(z) \cap N_Y(z) = \{0\}$ .



**Step 1:** Linearize *X*; intersect with *Y*.

Recast as set intersection: find a point z = (u, 0) where

 $X = \operatorname{gph} \Phi$  intersects  $Y = \mathbf{R}^n \times \{0\}.$ 

Assume transversality:  $N_X(z) \cap N_Y(z) = \{0\}$ .



**Step 1:** Linearize *X*; intersect with *Y*.

**Step 2:** Restore to *X* via a Lipschitz map fixing *z*.

Recast as set intersection: find a point z = (u, 0) where

 $X = \operatorname{gph} \Phi$  intersects  $Y = \mathbf{R}^n \times \{0\}.$ 

Assume transversality:  $N_X(z) \cap N_Y(z) = \{0\}$ .



**Step 1:** Linearize *X*; intersect with *Y*.

Linearize around  $v \in \Phi(u)$ ; solve for u'

**Step 2:** Restore to *X* via a Lipschitz map fixing *z*.

Recast as set intersection: find a point z = (u, 0) where

 $X = \operatorname{gph} \Phi$  intersects  $Y = \mathbf{R}^n \times \{0\}.$ 

Assume transversality:  $N_X(z) \cap N_Y(z) = \{0\}$ .





(L-Wylie '20)

**Step 1:** Linearize *X*; intersect with *Y*.

Linearize around  $v \in \Phi(u)$ ; solve for u'

**Step 2:** Restore to *X* via a Lipschitz map fixing *z*.

 $u^+ = \operatorname{Proj}_{\mathcal{M}}(u'); \quad v^+ = \operatorname{Proj}_{\Phi(u^+)}(0)$ 

# Fast black box nonsmooth optimization

Newtonian methods for partly smooth optimization  $0 \in \partial f(x)$  are interesting, but typically need structural knowledge of  $\partial f$ .

# Fast black box nonsmooth optimization

Newtonian methods for partly smooth optimization  $0 \in \partial f(x)$  are interesting, but typically need structural knowledge of  $\partial f$ .

Classical special case: sequential quadratic programming.

# Fast black box nonsmooth optimization

Newtonian methods for partly smooth optimization  $0 \in \partial f(x)$  are interesting, but typically need structural knowledge of  $\partial f$ .

Classical special case: sequential quadratic programming.

More generally, semismooth Newton methods: Klatte-Kummer '02, Facchinei-Pang '03, Izmailov-Solodov '14, Gfrerer-Outrata '19.

Newtonian methods for partly smooth optimization  $0 \in \partial f(x)$  are interesting, but typically need structural knowledge of  $\partial f$ .

Classical special case: sequential quadratic programming.

More generally, semismooth Newton methods: Klatte-Kummer '02, Facchinei-Pang '03, Izmailov-Solodov '14, Gfrerer-Outrata '19.

With just an oracle for **linear** approximations to convex f at input points, bundle methods are appealing (Sagastizábal '18 ICM).

Newtonian methods for partly smooth optimization  $0 \in \partial f(x)$  are interesting, but typically need structural knowledge of  $\partial f$ .

Classical special case: sequential quadratic programming.

More generally, semismooth Newton methods: Klatte-Kummer '02, Facchinei-Pang '03, Izmailov-Solodov '14, Gfrerer-Outrata '19.

With just an oracle for **linear** approximations to convex f at input points, bundle methods are appealing (Sagastizábal '18 ICM).

- "Null" steps enhance a cutting plane model.
- "Serious" steps sufficiently decrease the objective
- Partial smoothness ("VU") can accelerate the serious steps.

Newtonian methods for partly smooth optimization  $0 \in \partial f(x)$  are interesting, but typically need structural knowledge of  $\partial f$ .

Classical special case: sequential quadratic programming.

More generally, semismooth Newton methods: Klatte-Kummer '02, Facchinei-Pang '03, Izmailov-Solodov '14, Gfrerer-Outrata '19.

With just an oracle for **linear** approximations to convex f at input points, bundle methods are appealing (Sagastizábal '18 ICM).

- "Null" steps enhance a cutting plane model.
- "Serious" steps sufficiently decrease the objective
- Partial smoothness ("VU") can accelerate the serious steps.

But can we reduce oracle calls using quadratic approximations?

Convex  $f : \mathbf{R}^n \to \mathbf{R}$  are twice differentiable off a negligible set N.

Convex  $f : \mathbb{R}^n \to \mathbb{R}$  are twice differentiable off a negligible set N. Black-box methods (bundle, BFGS) typically never encounter N.

Convex  $f : \mathbb{R}^n \to \mathbb{R}$  are twice differentiable off a negligible set N. Black-box methods (bundle, BFGS) typically never encounter N. What if an oracle returns  $f(x), \nabla f(x), \nabla^2 f(x)$  for  $x \notin N$ ?

Convex  $f : \mathbb{R}^n \to \mathbb{R}$  are twice differentiable off a negligible set N. Black-box methods (bundle, BFGS) typically never encounter N. What if an oracle returns  $f(x), \nabla f(x), \nabla^2 f(x)$  for  $x \notin N$ ?

Aim: find a bundle S of k reference points, with small diameter

$$\max\{\|s-s'\|:s,s'\in S\}$$

and small optimality measure

$$dist(0, conv(\nabla f(S))).$$

Convex  $f : \mathbb{R}^n \to \mathbb{R}$  are twice differentiable off a negligible set N. Black-box methods (bundle, BFGS) typically never encounter N. What if an oracle returns  $f(x), \nabla f(x), \nabla^2 f(x)$  for  $x \notin N$ ?

Aim: find a bundle S of k reference points, with small diameter

$$\max\{\|s-s'\|:s,s'\in S\}$$

and small optimality measure

$$dist(0, conv(\nabla f(S))).$$

Intuition: if the bundle size k is large enough,

$$\lim_{S\to\{x\}}\operatorname{conv}(\nabla f(S)) = \partial f(x).$$

# (L-Wylie '19)

(L-Wylie '19)

For each of the k current reference points  $s \in S$ , use the oracle to form the linear and quadratic approximations

$$l_{s}(x) = f(s) + \nabla f(s)^{T}(x-s)$$
  

$$q_{s}(x) = l_{s}(x) + \frac{1}{2}(x-s)^{T} \nabla^{2} f(s)(x-s).$$

(L-Wylie '19)

For each of the k current reference points  $s \in S$ , use the oracle to form the linear and quadratic approximations

$$l_{s}(x) = f(s) + \nabla f(s)^{T}(x-s)$$
  
$$q_{s}(x) = l_{s}(x) + \frac{1}{2}(x-s)^{T} \nabla^{2} f(s)(x-s).$$

• Choose bundle weights  $\lambda_s$  solving

$$\min\Big\{\Big\|\sum_{s\in S}\lambda_s\nabla f(s)\Big\|:\lambda\geq 0,\ \sum_s\lambda_s=1\Big\}.$$

(L-Wylie '19)

For each of the k current reference points  $s \in S$ , use the oracle to form the linear and quadratic approximations

$$l_{s}(x) = f(s) + \nabla f(s)^{T}(x-s)$$
  
$$q_{s}(x) = l_{s}(x) + \frac{1}{2}(x-s)^{T} \nabla^{2} f(s)(x-s).$$

• Choose bundle weights  $\lambda_s$  solving

$$\min\Big\{\Big\|\sum_{s\in S}\lambda_s\nabla f(s)\Big\|:\lambda\geq 0,\ \sum_s\lambda_s=1\Big\}.$$

• Choose a new reference point x solving

$$\min\Big\{\sum_s\lambda_sq_s(x):l_s(s)\text{ equal for all }s\in S\Big\}.$$

(L-Wylie '19)

For each of the k current reference points  $s \in S$ , use the oracle to form the linear and quadratic approximations

$$l_{s}(x) = f(s) + \nabla f(s)^{T}(x-s)$$
  
$$q_{s}(x) = l_{s}(x) + \frac{1}{2}(x-s)^{T} \nabla^{2} f(s)(x-s).$$

• Choose bundle weights  $\lambda_s$  solving

$$\min\Big\{\Big\|\sum_{s\in S}\lambda_s\nabla f(s)\Big\|:\lambda\geq 0,\ \sum_s\lambda_s=1\Big\}.$$

• Choose a new reference point x solving

$$\min\Big\{\sum_{s}\lambda_{s}q_{s}(x):I_{s}(s)\text{ equal for all }s\in S\Big\}.$$

• Replace  $s \in S$  minimizing  $\|\nabla f(s) - \nabla f(x)\|$  with x.
Given a current bundle of reference points  $s \in S$ , model

$$f(x) \approx \max_{s \in S} f_s(x)$$
:

unknown smooth component  $f_s$  matches f to 2nd order around s.

Given a current bundle of reference points  $s \in S$ , model

$$f(x) \approx \max_{s \in S} f_s(x)$$
:

unknown smooth component  $f_s$  matches f to 2nd order around s. Optimize model via sequential quadratic programming steps on

$$\left\{egin{array}{ll} {
m minimize} & t \ {
m subject to} & f_s(x) \leq t \ {
m (} s \in S{
m )}. \end{array}
ight.$$

Given a current bundle of reference points  $s \in S$ , model

$$f(x) \approx \max_{s \in S} f_s(x)$$
:

unknown smooth component  $f_s$  matches f to 2nd order around s.

Optimize model via sequential quadratic programming steps on

$$\left\{egin{array}{ll} {
m minimize} & t \ {
m subject to} & f_s(x) \leq t \ {
m (} s \in S{
m )}. \end{array}
ight.$$

• Estimate multipliers via least squares.

Given a current bundle of reference points  $s \in S$ , model

$$f(x) \approx \max_{s \in S} f_s(x)$$
:

unknown smooth component  $f_s$  matches f to 2nd order around s.

Optimize model via sequential quadratic programming steps on

$$\left\{egin{array}{ll} {
m minimize} & t \ {
m subject to} & f_s(x) \leq t \ {
m (} s \in S{
m )}. \end{array}
ight.$$

- Estimate multipliers via least squares.
- Minimize the quadratic approximation of the Lagrangian...

Given a current bundle of reference points  $s \in S$ , model

$$f(x) \approx \max_{s \in S} f_s(x)$$
:

unknown smooth component  $f_s$  matches f to 2nd order around s.

Optimize model via sequential quadratic programming steps on

$$\left\{egin{array}{ll} {
m minimize} & t \ {
m subject to} & f_s(x) \leq t \ {
m (} s \in S{
m )}. \end{array}
ight.$$

- Estimate multipliers via least squares.
- Minimize the quadratic approximation of the Lagrangian...
- ... subject to the linearized constraints (assumed all active).

Given a current bundle of reference points  $s \in S$ , model

$$f(x) \approx \max_{s \in S} f_s(x)$$
 :

unknown smooth component  $f_s$  matches f to 2nd order around s.

Optimize model via sequential quadratic programming steps on

$$\left\{egin{array}{ll} {
m minimize} & t \ {
m subject to} & f_s(x) \leq t \ {
m (} s \in S{
m )}. \end{array}
ight.$$

- Estimate multipliers via least squares.
- Minimize the quadratic approximation of the Lagrangian...
- ... subject to the linearized constraints (assumed all active).

New point x improves model's most closely matching component.

**Theorem** If *f* decomposes as max function of degree *k*,

$$f(x) = \max_{i=1,\ldots,k} f_i(x),$$

where each component  $f_i$  is smooth

**Theorem** If *f* decomposes as max function of degree *k*,

$$f(x) = \max_{i=1,\ldots,k} f_i(x),$$

where each component  $f_i$  is smooth and strongly convex

**Theorem** If *f* decomposes as max function of degree *k*,

$$f(x) = \max_{i=1,\ldots,k} f_i(x),$$

where each component  $f_i$  is smooth and strongly convex around a strictly active critical point  $\bar{x}$ ,

**Theorem** If *f* decomposes as max function of degree *k*,

$$f(x) = \max_{i=1,\ldots,k} f_i(x),$$

where each component  $f_i$  is smooth and strongly convex around a strictly active critical point  $\bar{x}$ , and initial  $S = \{s_1, \dots, s_k\}$  is a full bundle near  $\bar{x}$ , meaning  $f_i(s_i) > f_i(s_i)$  whenever  $i \neq j$ ,

**Theorem** If f decomposes as max function of degree k,

$$f(x) = \max_{i=1,\ldots,k} f_i(x),$$

where each component  $f_i$  is smooth and strongly convex around a strictly active critical point  $\bar{x}$ , and initial  $S = \{s_1, \dots, s_k\}$  is a full bundle near  $\bar{x}$ , meaning  $f_i(s_i) > f_i(s_i)$  whenever  $i \neq j$ ,

then k-bundle Newton converges k-step quadratically to  $\bar{x}$ .

**Theorem** If f decomposes as max function of degree k,

$$f(x) = \max_{i=1,\ldots,k} f_i(x),$$

where each component  $f_i$  is smooth and strongly convex around a strictly active critical point  $\bar{x}$ , and initial  $S = \{s_1, \dots, s_k\}$  is a full bundle near  $\bar{x}$ , meaning  $f_i(s_i) > f_i(s_i)$  whenever  $i \neq j$ ,

then k-bundle Newton converges k-step quadratically to  $\bar{x}$ .

**Note**: The required bundle size k and the partly smooth geometry of the active manifold M are related:

$$k + \dim \mathcal{M} = n + 1.$$

**Theorem** If *f* decomposes as max function of degree *k*,

$$f(x) = \max_{i=1,\ldots,k} f_i(x),$$

where each component  $f_i$  is smooth and strongly convex around a strictly active critical point  $\bar{x}$ , and initial  $S = \{s_1, \dots, s_k\}$  is a full bundle near  $\bar{x}$ , meaning  $f_i(s_i) > f_i(s_i)$  whenever  $i \neq j$ ,

then k-bundle Newton converges k-step quadratically to  $\bar{x}$ .

**Note**: The required bundle size k and the partly smooth geometry of the active manifold  $\mathcal{M}$  are related:

$$k + \dim \mathcal{M} = n + 1.$$

The classical Newton method has k = 1.

$$f(x,y) = 2x^2 + y^2 + |x^2 - y|$$

$$f(x,y) = 2x^2 + y^2 + |x^2 - y|$$

Objective value against oracle calls.

$$f(x,y) = 2x^2 + y^2 + |x^2 - y|$$

Objective value against oracle calls.

• Naive proximal bundle \*

$$f(x,y) = 2x^2 + y^2 + |x^2 - y|$$

Objective value against oracle calls.

- Naive proximal bundle \*
- Nonsmooth BFGS  $\circ$

$$f(x,y) = 2x^2 + y^2 + |x^2 - y|$$

Objective value against oracle calls.

- Naive proximal bundle \*
- Nonsmooth BFGS  $\circ$
- 2-bundle Newton

$$f(x,y) = 2x^2 + y^2 + |x^2 - y|$$

Objective value against oracle calls.

- Naive proximal bundle \*
- Nonsmooth BFGS  $\circ$
- 2-bundle Newton (initiated from

proximal bundle)

$$f(x,y) = 2x^2 + y^2 + |x^2 - y|$$

Objective value against oracle calls.

- Naive proximal bundle \*
- Nonsmooth BFGS  $\circ$
- 2-bundle Newton (initiated from proximal bundle)
  - objective **X**

$$f(x,y) = 2x^2 + y^2 + |x^2 - y|$$

Objective value against oracle calls.

- Naive proximal bundle \*
- Nonsmooth BFGS  $\circ$
- 2-bundle Newton

(initiated from proximal bundle)

- objective **X**
- optimality measure  $\triangle$

$$f(x,y) = 2x^2 + y^2 + |x^2 - y|$$

Objective value against oracle calls.

- Naive proximal bundle \*
- Nonsmooth BFGS  $\circ$
- 2-bundle Newton

(initiated from proximal bundle)

- objective **X**
- optimality measure  $\triangle$
- bundle diameter  $\diamond$

$$f(x,y) = 2x^2 + y^2 + |x^2 - y|$$



$$f(s, t, u, v) = \lambda_{\max} \begin{bmatrix} 0 & s & t \\ s & 1+u & v \\ t & v & 1-u \end{bmatrix}$$

$$f(s,t,u,v) = \lambda_{\max} \begin{bmatrix} 0 & s & t \\ s & 1+u & v \\ t & v & 1-u \end{bmatrix} - 1$$

Not a max function around its minimizer, zero.

$$f(s,t,u,v) = \lambda_{\max} \begin{bmatrix} 0 & s & t \\ s & 1+u & v \\ t & v & 1-u \end{bmatrix} - 1$$

Not a max function around its minimizer, zero. But...

... partly smooth relative to a 2-dimensional manifold.

$$f(s,t,u,v) = \lambda_{\max} \begin{bmatrix} 0 & s & t \\ s & 1+u & v \\ t & v & 1-u \end{bmatrix} - 1$$

Not a max function around its minimizer, zero. But...

... partly smooth relative to a 2-dimensional manifold.

- Proximal bundle \*
- BFGS •
- 3-bundle Newton
  - $\bullet~$  objective  ${\boldsymbol X}$
  - optimality  $\triangle$
  - diameter  $\diamond$

$$f(s, t, u, v) = \lambda_{\max} \begin{bmatrix} 0 & s & t \\ s & 1+u & v \\ t & v & 1-u \end{bmatrix} - 1$$

Not a max function around its minimizer, zero. But...

... partly smooth relative to a 2-dimensional manifold.

- Proximal bundle \*
- BFGS •
- 3-bundle Newton
  - objective **X**
  - optimality  $\triangle$
  - diameter  $\diamond$



Consider a maximum of smooth strongly convex components,

$$f(x) = \max_{i=1,\ldots,k} f_i(x),$$

with a strictly active critical point  $\bar{x}$ . on the active manifold  $\mathcal{M}$ .

Consider a maximum of smooth strongly convex components,

$$f(x) = \max_{i=1,\ldots,k} f_i(x),$$

with a strictly active critical point  $\bar{x}$ . on the active manifold  $\mathcal{M}$ .

Near x on the active manifold  $\mathcal{M}$ , full bundles  $S = \{s_1, \ldots, s_k\}$ (so  $f_i(s_i) > f_j(s_i)$  for  $i \neq j$ ) approximate the subdifferential:

 $\partial f(x) \approx \operatorname{conv} \nabla f(S),$ 

Consider a maximum of smooth strongly convex components,

$$f(x) = \max_{i=1,\ldots,k} f_i(x),$$

with a strictly active critical point  $\bar{x}$ . on the active manifold  $\mathcal{M}$ .

Near x on the active manifold  $\mathcal{M}$ , full bundles  $S = \{s_1, \ldots, s_k\}$ (so  $f_i(s_i) > f_j(s_i)$  for  $i \neq j$ ) approximate the subdifferential:

 $\partial f(x) \approx \operatorname{conv} \nabla f(S),$ 

and the Hessians  $\nabla^2 f(s_i)$  predict curvature on  $\mathcal{M}$ .

Consider a maximum of smooth strongly convex components,

$$f(x) = \max_{i=1,\ldots,k} f_i(x),$$

with a strictly active critical point  $\bar{x}$ . on the active manifold  $\mathcal{M}$ .

Near x on the active manifold  $\mathcal{M}$ , full bundles  $S = \{s_1, \ldots, s_k\}$ (so  $f_i(s_i) > f_j(s_i)$  for  $i \neq j$ ) approximate the subdifferential:

$$\partial f(x) \approx \operatorname{conv} \nabla f(S),$$

and the Hessians  $\nabla^2 f(s_i)$  predict curvature on  $\mathcal{M}$ .

Partly smooth geometry then ensures  $\hat{x} - \bar{x} = O(|\bar{x} - S|^2)$ .

Consider a maximum of smooth strongly convex components,

$$f(x) = \max_{i=1,\ldots,k} f_i(x),$$

with a strictly active critical point  $\bar{x}$ . on the active manifold  $\mathcal{M}$ .

Near x on the active manifold  $\mathcal{M}$ , full bundles  $S = \{s_1, \ldots, s_k\}$ (so  $f_i(s_i) > f_j(s_i)$  for  $i \neq j$ ) approximate the subdifferential:

$$\partial f(x) \approx \operatorname{conv} \nabla f(S),$$

and the Hessians  $\nabla^2 f(s_i)$  predict curvature on  $\mathcal{M}$ .

Partly smooth geometry then ensures  $\hat{x} - \bar{x} = O(|\bar{x} - S|^2)$ .

Updating  $s_i \leftarrow \hat{x}$  keeps the bundle full, because  $\bar{x}$  is strictly active.

Consider a maximum of smooth strongly convex components,

$$f(x) = \max_{i=1,\ldots,k} f_i(x),$$

with a strictly active critical point  $\bar{x}$ . on the active manifold  $\mathcal{M}$ .

Near x on the active manifold  $\mathcal{M}$ , full bundles  $S = \{s_1, \ldots, s_k\}$ (so  $f_i(s_i) > f_j(s_i)$  for  $i \neq j$ ) approximate the subdifferential:

$$\partial f(x) \approx \operatorname{conv} \nabla f(S),$$

and the Hessians  $\nabla^2 f(s_i)$  predict curvature on  $\mathcal{M}$ .

Partly smooth geometry then ensures  $\hat{x} - \bar{x} = O(|\bar{x} - S|^2)$ .

Updating  $s_i \leftarrow \hat{x}$  keeps the bundle full, because  $\bar{x}$  is strictly active.

Each reference point  $s_i$  updates within k steps, by strong convexity.

## Finding an initial bundle

Black-box methods for finding a minimizer  $\bar{x}$  for nonsmooth f, like

- bundle methods (Lemaréchal, Wolfe '70's)
- BFGS (L-Overton '13)
- gradient sampling (Burke-L-Overton '05)

## Finding an initial bundle

Black-box methods for finding a minimizer  $\bar{x}$  for nonsmooth f, like

- bundle methods (Lemaréchal, Wolfe '70's)
- BFGS (L-Overton '13)
- gradient sampling (Burke-L-Overton '05)

asymptotically generate subdifferential approximations:

$$\partial f(\bar{x}) \approx \operatorname{conv}(\nabla f(\Omega))$$

for sets  $\Omega$  of points near  $\bar{x}$ .
# Finding an initial bundle

Black-box methods for finding a minimizer  $\bar{x}$  for nonsmooth f, like

- bundle methods (Lemaréchal, Wolfe '70's)
- BFGS (L-Overton '13)
- gradient sampling (Burke-L-Overton '05)

asymptotically generate subdifferential approximations:

$$\partial f(\bar{x}) \approx \operatorname{conv}(\nabla f(\Omega))$$

for sets  $\Omega$  of points near  $\bar{x}$ . So, we could choose

$$k \;=\; \dim \Bigl( {
m affine \; span} ig( 
abla f(\Omega) ig) \Bigr) \qquad ({
m numerically})$$

# Finding an initial bundle

Black-box methods for finding a minimizer  $\bar{x}$  for nonsmooth f, like

- bundle methods (Lemaréchal, Wolfe '70's)
- BFGS (L-Overton '13)
- gradient sampling (Burke-L-Overton '05)

asymptotically generate subdifferential approximations:

$$\partial f(\bar{x}) \approx \operatorname{conv}(\nabla f(\Omega))$$

for sets  $\Omega$  of points near  $\bar{x}$ . So, we could choose

$$k \;=\; \dim \Bigl( {
m affine \; span} ig( 
abla f(\Omega) ig) \Bigr) \qquad ({
m numerically})$$

and initial  $S \subset \Omega$  of size k with  $\nabla f(S)$  affinely independent.

For random 25-by-25 symmetric matrices, minimize  $\lambda_{\max}$  for

$$A(x) = A_0 + x_1A_1 + x_2A_2 + \ldots + x_{50}A_{50}.$$

For random 25-by-25 symmetric matrices, minimize  $\lambda_{\max}$  for

$$A(x) = A_0 + x_1A_1 + x_2A_2 + \ldots + x_{50}A_{50}.$$

Active manifold, where  $\lambda_{\max}(A(x))$  has multiplicity 6, has dim 30.

For random 25-by-25 symmetric matrices, minimize  $\lambda_{\max}$  for

$$A(x) = A_0 + x_1A_1 + x_2A_2 + \ldots + x_{50}A_{50}.$$

Active manifold, where  $\lambda_{\max}(A(x))$  has multiplicity 6, has dim 30.





 $f - \min f$ 

For random 25-by-25 symmetric matrices, minimize  $\lambda_{\max}$  for

$$A(x) = A_0 + x_1A_1 + x_2A_2 + \ldots + x_{50}A_{50}.$$

Active manifold, where  $\lambda_{\max}(A(x))$  has multiplicity 6, has dim 30.

 $\lambda_1(A(x)) - \lambda_6(A(x))$ 



# Extensions...

• Avoiding Hessians...

- Avoiding Hessians...
  - ... using automatic differentiation

- Avoiding Hessians...
  - $\bullet$   $\ldots$  using automatic differentiation  $\checkmark$

- Avoiding Hessians...
  - $\bullet$   $\ldots$  using automatic differentiation  $\checkmark$
  - ... with a linearly convergent first-order analogue

- Avoiding Hessians...
  - $\bullet$   $\ldots$  using automatic differentiation  $\checkmark$
  - ... with a linearly convergent first-order analogue  $\checkmark?$

- Avoiding Hessians...
  - $\bullet$   $\ldots$  using automatic differentiation  $\checkmark$
  - ... with a linearly convergent first-order analogue  $\checkmark$ ?
- Extending the local convergence analysis...

- Avoiding Hessians...
  - $\bullet$   $\ldots$  using automatic differentiation  $\checkmark$
  - ... with a linearly convergent first-order analogue  $\checkmark$ ?
- Extending the local convergence analysis...
  - $\bullet$   $\ldots$  to nonconvex max functions

- Avoiding Hessians...
  - $\bullet$   $\ldots$  using automatic differentiation  $\checkmark$
  - ... with a linearly convergent first-order analogue  $\checkmark$ ?
- Extending the local convergence analysis...
  - $\bullet$  ...to nonconvex max functions  $\checkmark$

- Avoiding Hessians...
  - $\bullet$   $\ldots$  using automatic differentiation  $\checkmark$
  - ... with a linearly convergent first-order analogue  $\checkmark$ ?
- Extending the local convergence analysis...
  - $\bullet$  ... to nonconvex max functions  $\checkmark$
  - $\bullet$   $\ldots$  to partly smooth functions

- Avoiding Hessians...
  - $\bullet$   $\ldots$  using automatic differentiation  $\checkmark$
  - ... with a linearly convergent first-order analogue  $\checkmark$ ?
- Extending the local convergence analysis...
  - $\bullet$  ... to nonconvex max functions  $\checkmark$
  - ... to partly smooth functions ??

- Avoiding Hessians...
  - $\bullet$   $\ldots$  using automatic differentiation  $\checkmark$
  - ... with a linearly convergent first-order analogue  $\checkmark$ ?
- Extending the local convergence analysis...
  - $\bullet$  ... to nonconvex max functions  $\checkmark$
  - ... to partly smooth functions ??
- Globalizing the algorithm

- Avoiding Hessians...
  - $\bullet$   $\ldots$  using automatic differentiation  $\checkmark$
  - ... with a linearly convergent first-order analogue  $\checkmark$ ?
- Extending the local convergence analysis...
  - $\bullet$  ... to nonconvex max functions  $\checkmark$
  - ... to partly smooth functions ??
- Globalizing the algorithm ??

Partial smoothness is a simple differential-geometric idea that captures the generic interplay between smooth and nonsmooth geometry in concrete variational problems, illuminating the analysis and design of algorithms. L and M.L. Overton, "Partial smoothness of the numerical radius at matrices whose fields of values are disks", *SIMAX* 2020.

X.Y. Han and L, "Disk matrices and the proximal mapping for the numerical radius", arXiv:2004.14542

L and J. Liang, "Partial smoothness and constant rank", arXiv:1807.03134.

L and C.J.S. Wylie, "Active-set Newton methods and partial smoothness", *MOR* 2020.

L and C.J.S. Wylie, "A simple Newton method for local nonsmooth optimization", arxiv:1907.11742.