

Combining Newton and proximal-gradient for nonsmooth optimization

Gilles Bareilles

LJK, Univ. Grenoble Alpes

CANUM
13 - 17 juin 2022

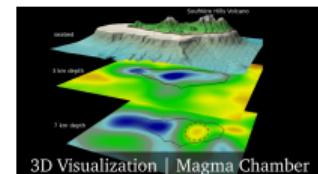
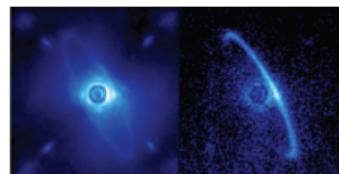
Context: linear inverse problems...

x_0 is the signal of interest, only accessible through measures y , via a (linear) model A and noise ξ :

$$y = Ax_0 + \xi$$

Examples

- ▶ **Image processing** diffraction of objective, low-res./damaged sensor
- ▶ **Medical imaging** computerized tomography, MRI, EEG
- ▶ **Seismic imaging** wave propagation
- ▶ **Machine learning / Statistics** regression



Context: linear inverse problems... are often ill posed

$$y = Ax_0 + \xi$$

Problem: not enough observations y to recover x_0 or y inconsistent.

Observation: we often know the solution to be sparse, piecewise constant, low-rank (matrix); we have **prior knowledge** on the **structure** of x_0 .

How to leverage it?

Context: linear inverse problems... are often ill posed

$$y = Ax_0 + \xi$$

Problem: not enough observations y to recover x_0 or y inconsistent.

Observation: we often know the solution to be sparse, piecewise constant, low-rank (matrix); we have **prior knowledge** on the **structure** of x_0 .

How to leverage it? Minimize the regularized empirical risk

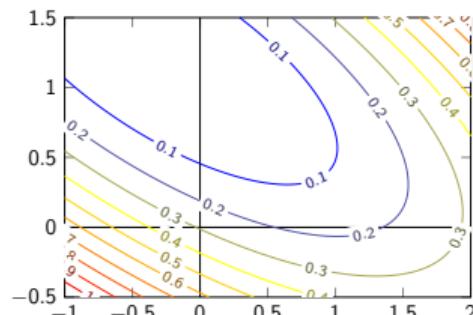
$$x^* \in \arg \min_x \underbrace{\sum_{i=1}^n \ell(A_i x, y_i)}_{\text{smooth empirical risk}} + \underbrace{\lambda r(x)}_{\text{nonsmooth regularizer}}$$

where

- ▶ ℓ is built from assumptions on noise ξ
- ▶ r is built to encode prior knowledge on structure of x_0

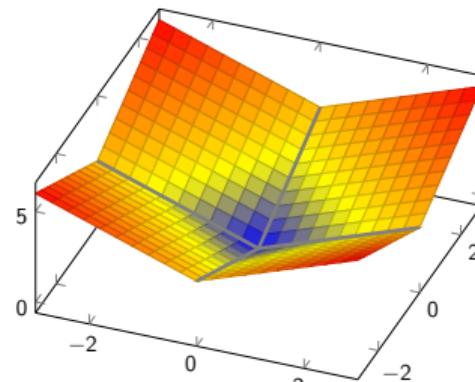
Nonsmooth regularizers and structure

▷ if x_0 has many zero entries, use $\|x\|_1 = \sum_{i=1}^n |x_i|$:



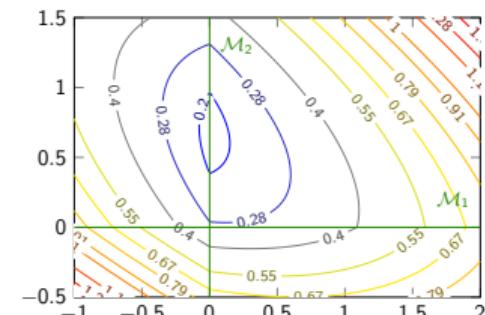
$$\sum_{i=1}^n (A_i x - y_i)^2$$

+



$$\lambda \|x\|_1$$

→



$$\sum_{i=1}^n (A_i x - y_i)^2 + \lambda \|x\|_1$$

LASSO

▷ if the matrix X_0 is low-rank, use $\|X\|_* = \sum_{i=1}^{\text{rank}(X)} \sigma_i(X)$

Observation: there are smooth subspaces M such that r is *smooth along* and *nonsmooth across*. These are **structure manifolds**.

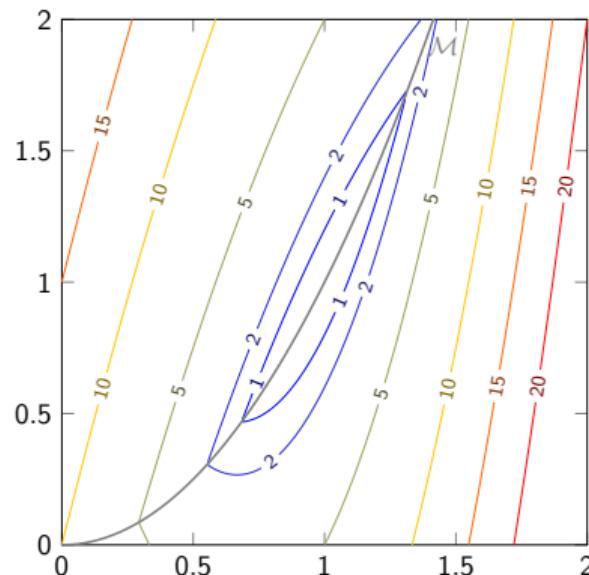
Composite problem, structure of minimizers

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

Finding a minimizer of F nonsmooth amounts to:

- ▶ finding the right structure manifold $\mathcal{M}^* \ni x^*$
e.g. the right sparsity pattern for ℓ_1
- ▶ minimizing F on that correct structure
smooth problem on a smooth manifold

But, we never know when a manifold is the optimal one.



Composite problem, structure of minimizers

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

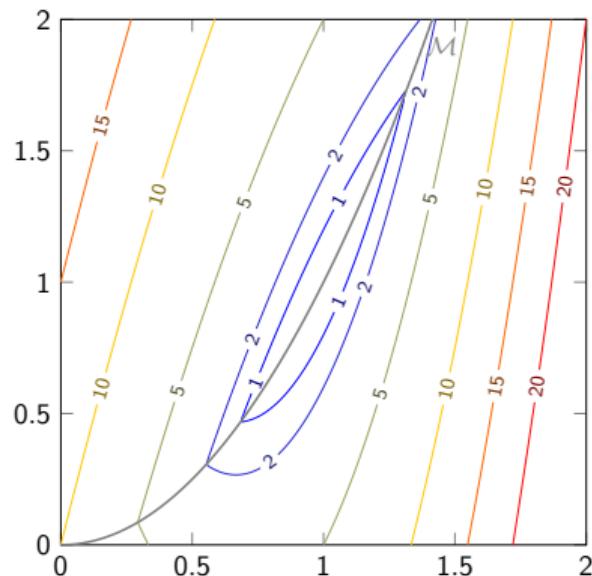
Finding a minimizer of F nonsmooth amounts to:

- ▶ finding the right structure manifold $\mathcal{M}^* \ni x^*$
e.g. the right sparsity pattern for ℓ_1
- ▶ minimizing F on that correct structure
smooth problem on a smooth manifold

But, we never know when a manifold is the optimal one.

Plan:

1. look at proximal gradient
2. look at Riemannian optimization
3. combine them



Proximal gradient: an optimization tool

- ▷ One proximal gradient step:

$$x_{k+1} = \mathbf{prox}_{\gamma g}(x_k - \gamma \nabla f(x_k)),$$

where

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

- ▷ Forward-backward interpretation:

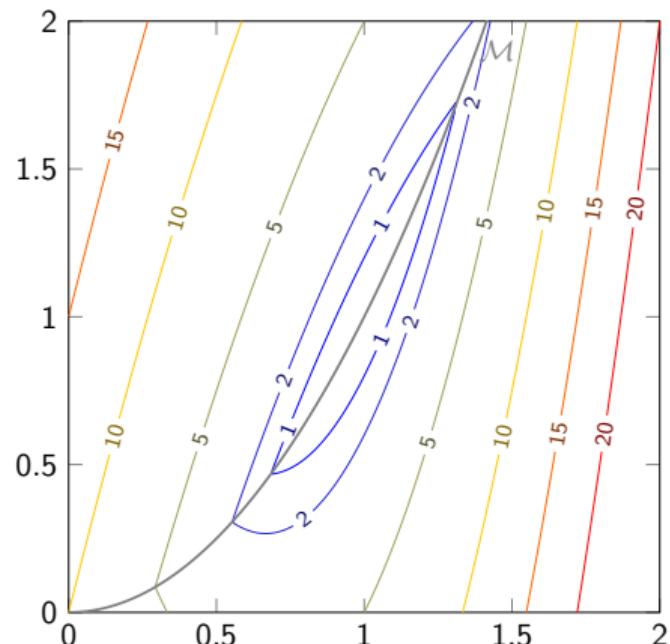
explicit step on f	$y = x - \gamma \nabla f(x)$
'implicit' step on g	$x_+ = \mathbf{prox}_{\gamma g}(y) \Leftrightarrow x_+ = y - \gamma s$ with $s \in \partial g(x_+)$

- ▷ If the step size γ is small enough $< 1/L$, where L is a Lipschitz constant for ∇f

- ✓ functional descent at every step hence converges to critical points
- ✗ “slow” convergence: $F(x_k) - F(x^*) \leq O(1/k)$ in convex case

Proximal gradient: a structure detector

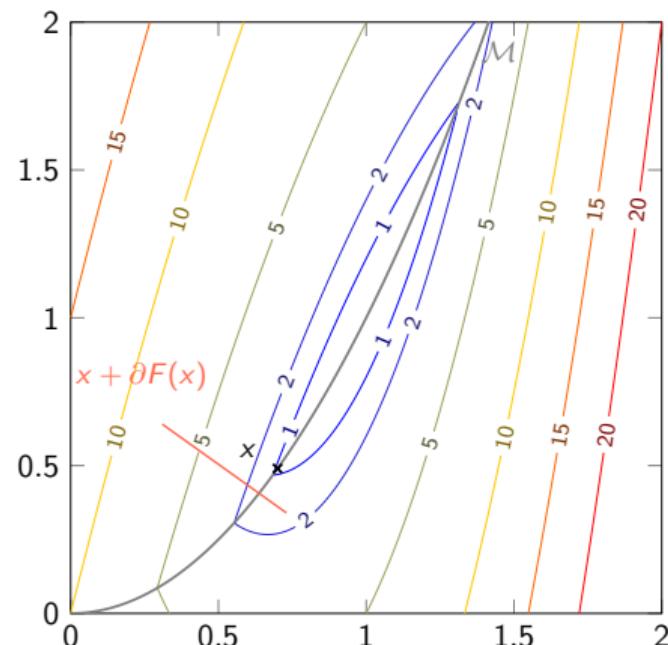
Near manifolds, the proximal gradient sends points to the manifold, smoothly:



$$g(x) = 10(1 - x_1)^2 + 5|x_1 - x_2|$$

Proximal gradient: a structure detector

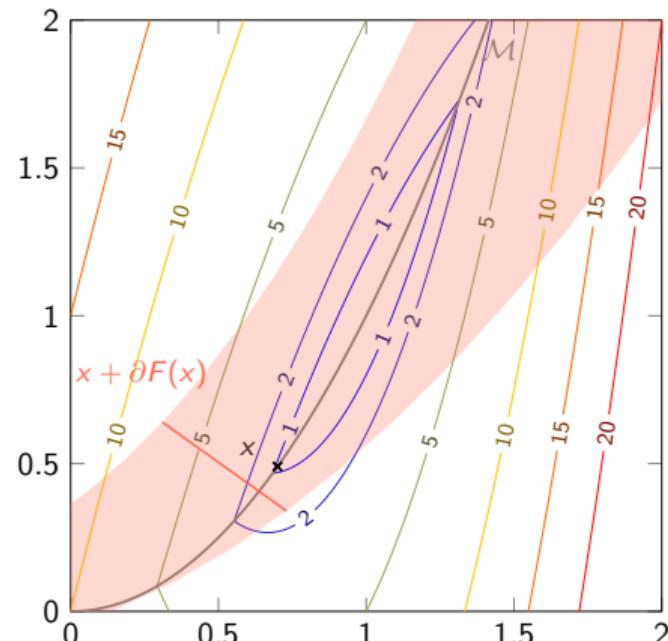
Near manifolds, the proximal gradient sends points to the manifold, smoothly:



$$g(x) = 10(1 - x_1)^2 + 5|x_1 - x_2|$$

Proximal gradient: a structure detector

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



$$g(x) = 10(1 - x_1)^2 + 5|x_1 - x_2|$$

Proximal gradient: a structure detector

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

Theorem (B., lutzeler, Malick, '20)

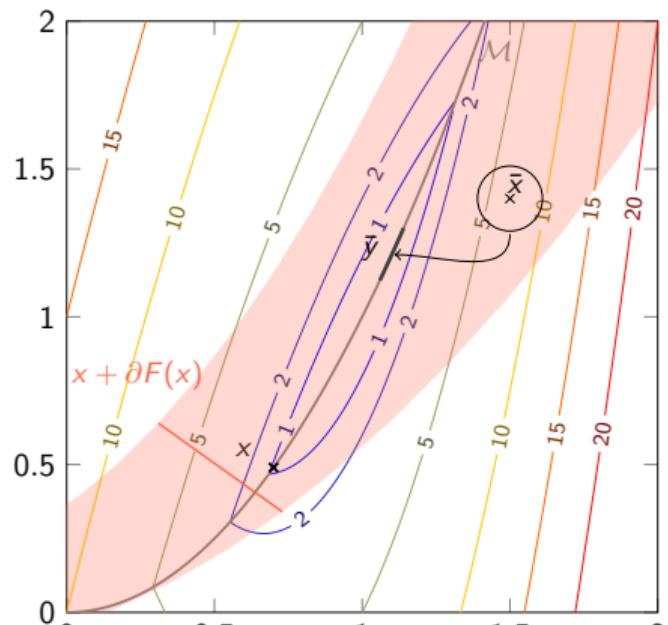
Take \bar{x} and $\bar{y} = \text{prox}_{\gamma g}(\bar{x} - \gamma \nabla f(\bar{x}))$ such that

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

then, on a neighborhood $\mathcal{N}_{\bar{x}}$ of \bar{x}

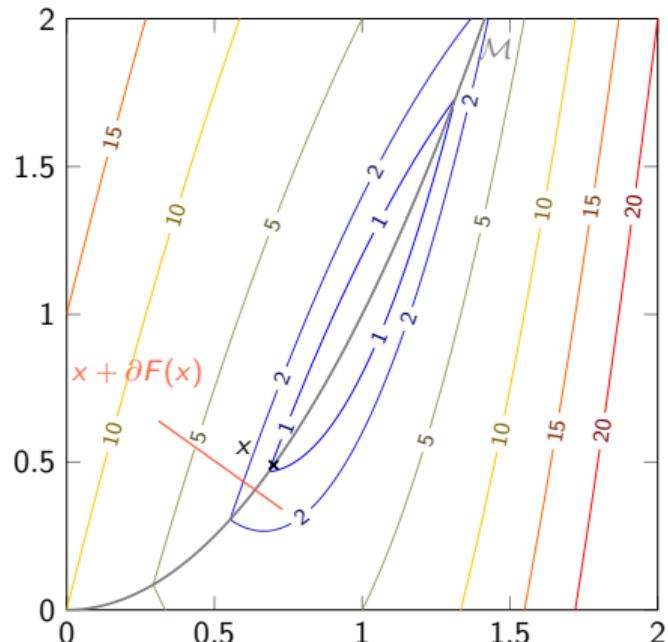
- ▶ the prox-grad operator is **\mathcal{M} -valued** and \mathcal{C}^1

→ How to move on \mathcal{M} ?



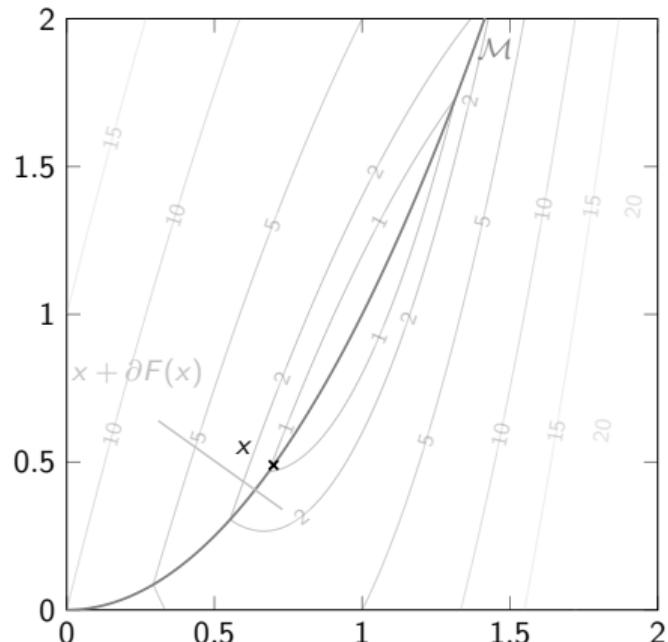
$$g(x) = 10(1-x_1)^2 + 5|x_1-x_2|$$

Manifold optimization, in a nutshell



$$\arg \min_x 10(1 - x_1)^2 + 5|x_1^2 - x_2|$$

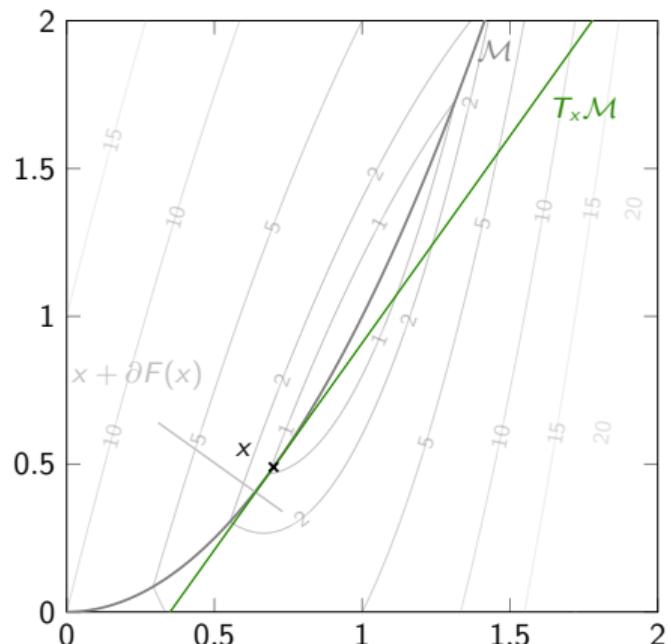
Manifold optimization, in a nutshell



$$\arg \min_x 10(1 - x_1)^2 + 5|x_1^2 - x_2|$$

Manifold optimization, in a nutshell

Elementary tools:

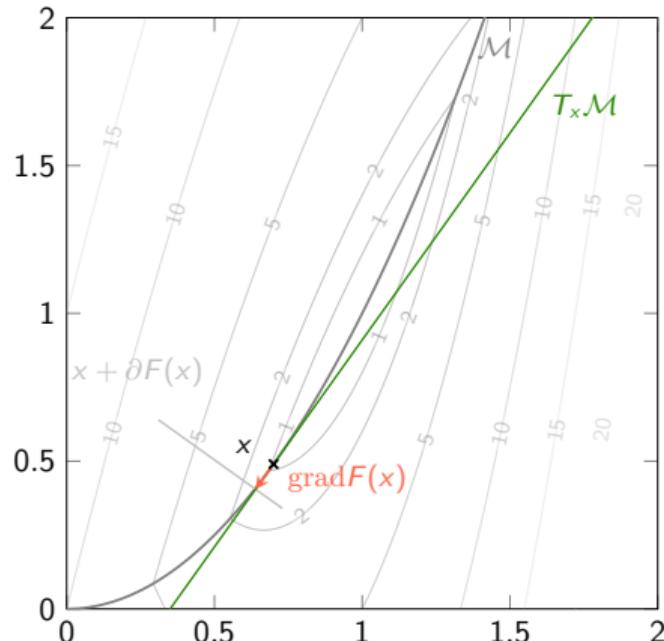


$$\arg \min_x 10(1 - x_1)^2 + 5|x_1^2 - x_2|$$

Manifold optimization, in a nutshell

Elementary tools:

- ▶ Riemannian gradient
- ▶ Riemannian Hessian

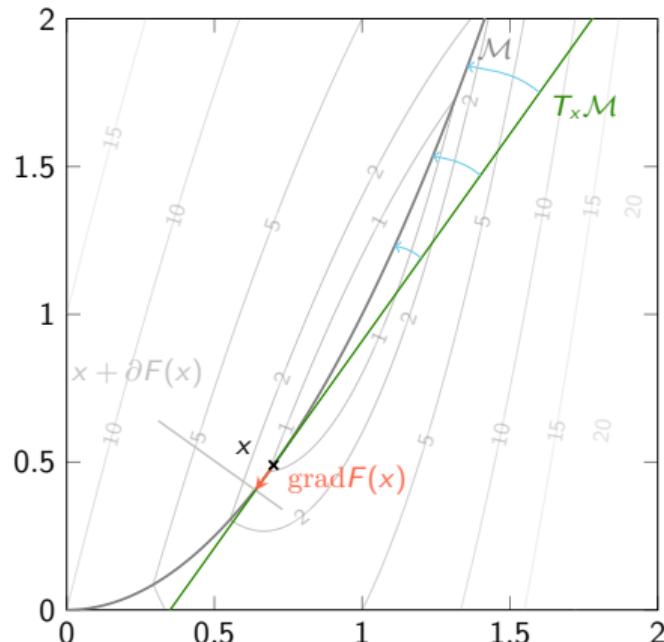


$$\arg \min_x 10(1 - x_1)^2 + 5|x_1^2 - x_2|$$

Manifold optimization, in a nutshell

Elementary tools:

- ▶ Riemannian gradient
- ▶ Riemannian Hessian
- ▶ retraction: $T_x\mathcal{M} \rightarrow \mathcal{M}$



$$\arg \min_x 10(1 - x_1)^2 + 5|x_1^2 - x_2|$$

Manifold optimization, in a nutshell

Elementary tools:

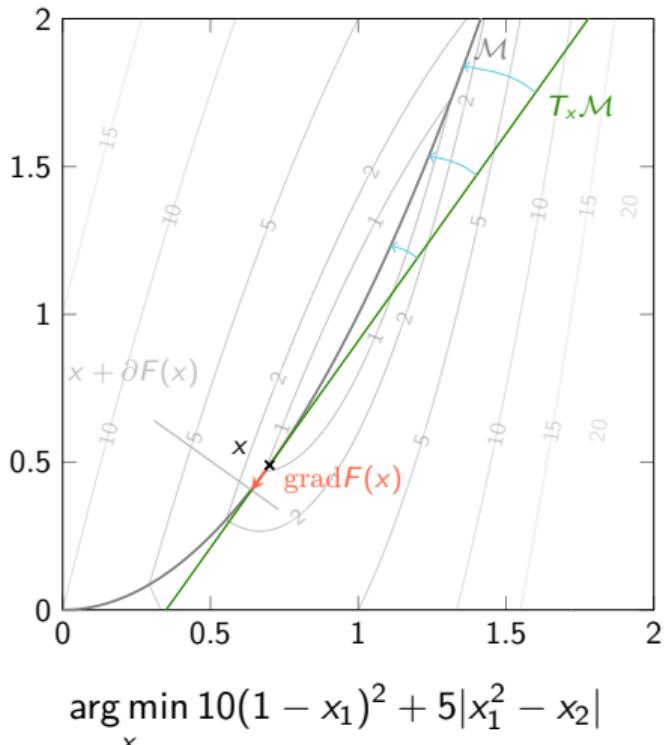
- ▶ Riemannian gradient
- ▶ Riemannian Hessian
- ▶ retraction: $T_x\mathcal{M} \rightarrow \mathcal{M}$

Typical Riemannian step “ManUp(x, \mathcal{M})”:

- ▶ find $d \in T_x\mathcal{M}$
- ▶ find a good step size $\alpha > 0$
- ▶ send αd to \mathcal{M} with a retraction

→ Many methods of smooth optim carry over manifold optimization.

◊ Boumal, N.: An introduction to optimization on smooth manifolds (2020).



Proposed algorithm

$\text{prox}_{\gamma g}(y_{k-1} - \gamma \nabla f(y_{k-1}))$ gives x_k and $\mathcal{M}_k \ni x_k$
 $y_k = \text{ManUp}(x_k, \mathcal{M}_k)$

Theorem (B., Iutzeler, Malick, '20)

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

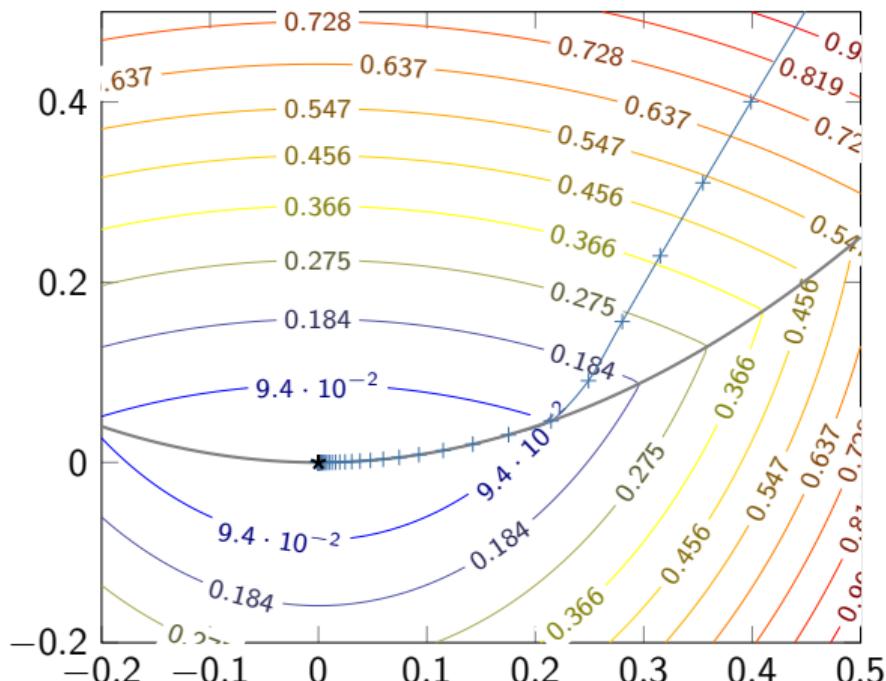
Take ManUp as a Riemannian Newton method and assume for a limit point x^* that

- ▶ 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**: $\|x_{k+1} - x^*\| \leq C \|x_k - x^*\|^2$

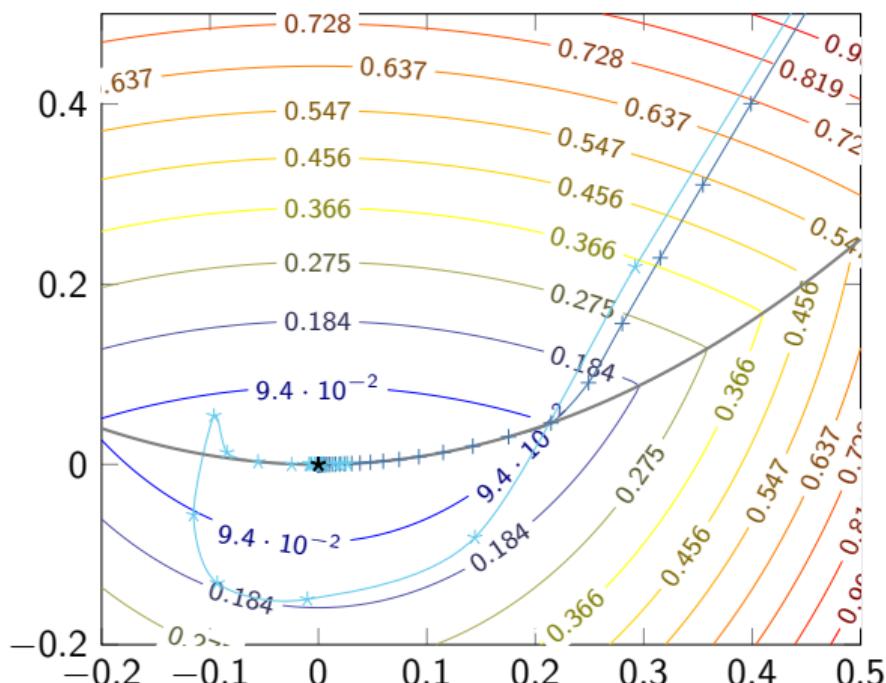
Illustrations



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

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

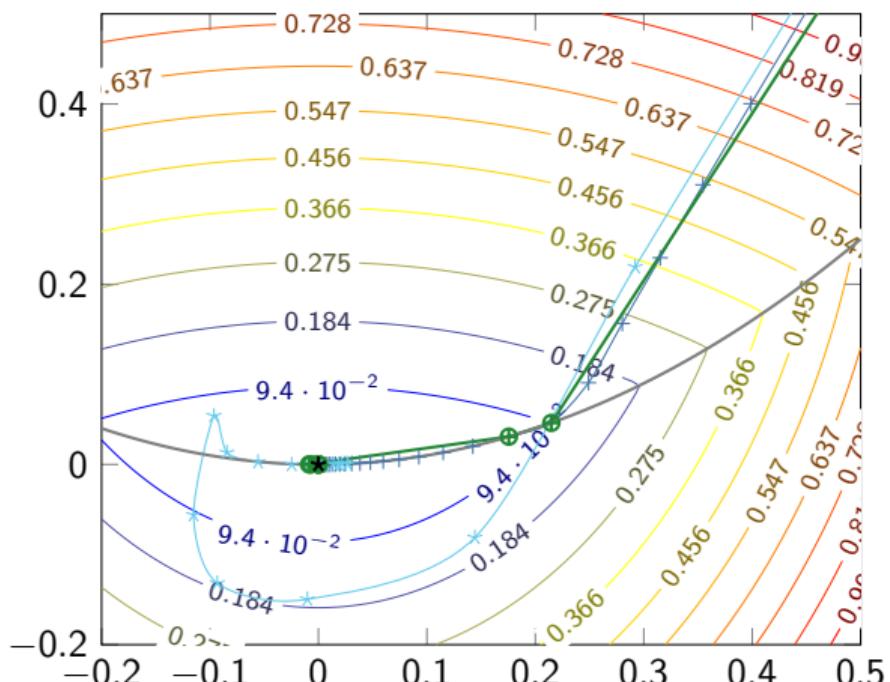
Illustrations



- Proximal Gradient
- Accel. Proximal Gradient
- Alt. Newton

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

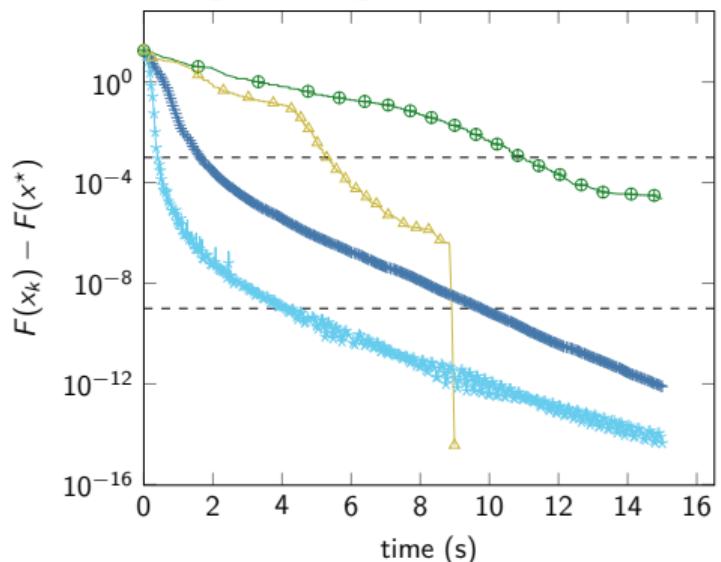
Illustrations



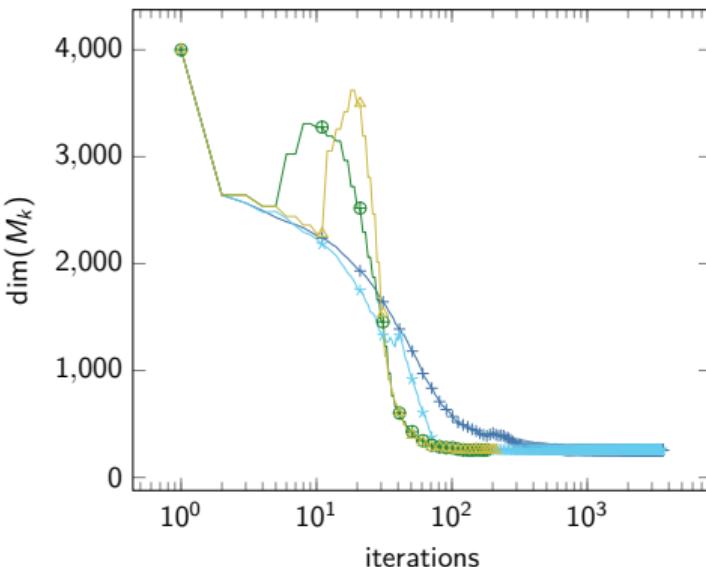
- +— Proximal Gradient
- *— Accel. Proximal Gradient
- ⊕— Alt. Newton

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

Illustrations: logistic regression



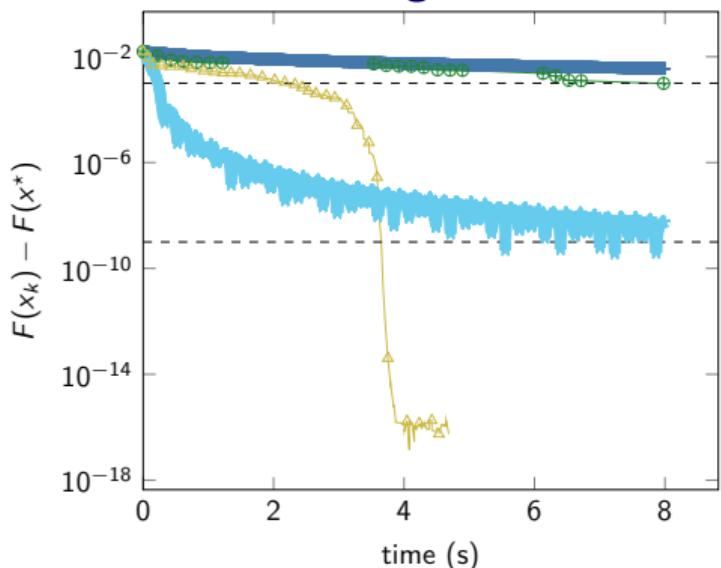
+	Proximal Gradient
*	Accel. Proximal Gradient
+	Alt. Newton
△	Alt. Truncated Newton



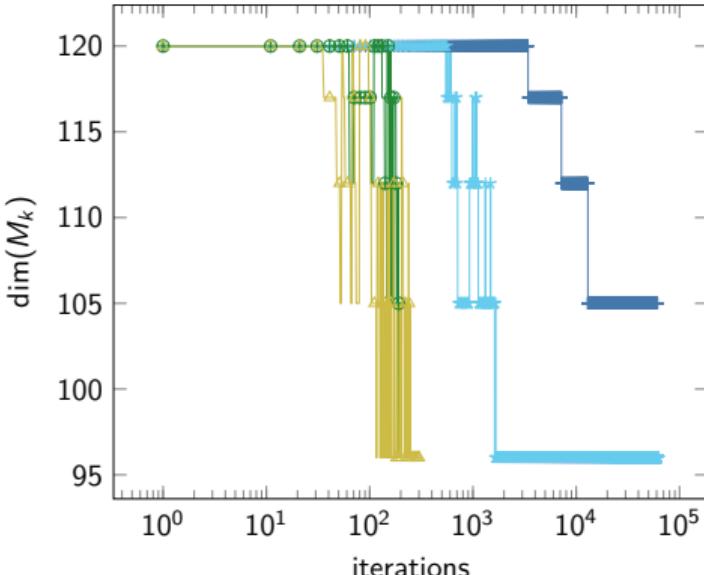
$$\min_{x \in \mathbb{R}^{4000}} \frac{1}{4000} \sum_{i=1}^{4000} \log(y_i \sigma(\langle A_i, x \rangle)) + 10^{-2} \|x\|_1,$$

The optimal manifold has dimension 249.

Illustrations: tracenorm regression



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



$$\min_{X \in \mathbb{R}^{10 \times 12}} \frac{1}{2} \sum_{i=1}^{60} (\langle A_i, X \rangle - y_i)^2 + 10^{-2} \|X\|_*,$$

The optimal manifold is $\text{rank}(X^*) = 6$.

Take-away messages

- ▶ Proximal methods identify structure in composite additive problems
- ▶ We proposed to leverage this structure with Newton-like Riemannian steps

B. & Iutzeler & Malick: *Newton acceleration on manifolds identified by proximal-gradient methods*, Mathematical Programming, 2nd revision.

→ <https://arxiv.org/abs/2012.12936>

→ <https://github.com/GillesBareilles/NewtonRiemannAccel-ProxGrad>

Work in progress and perspectives

- ▶ Smarter manifold updates
- ▶ Explore composite structures beyond addition $f + g$, e.g. $f + g + h(L \cdot)$ or composition $g \circ c$
- ▶ Certify that a structure is optimal

Thank you!