

Iteration complexity of optimization on smooth manifolds

Nicolas Boumal (Princeton University), joint work with P.-A. Absil (UCLouvain), Naman Agarwal (Google), Brian Bullins (Princeton), Coralia Cartis (Oxford)

Mostly harmless

Optimization on smooth manifolds is not that different from unconstrained optimization on a linear space:

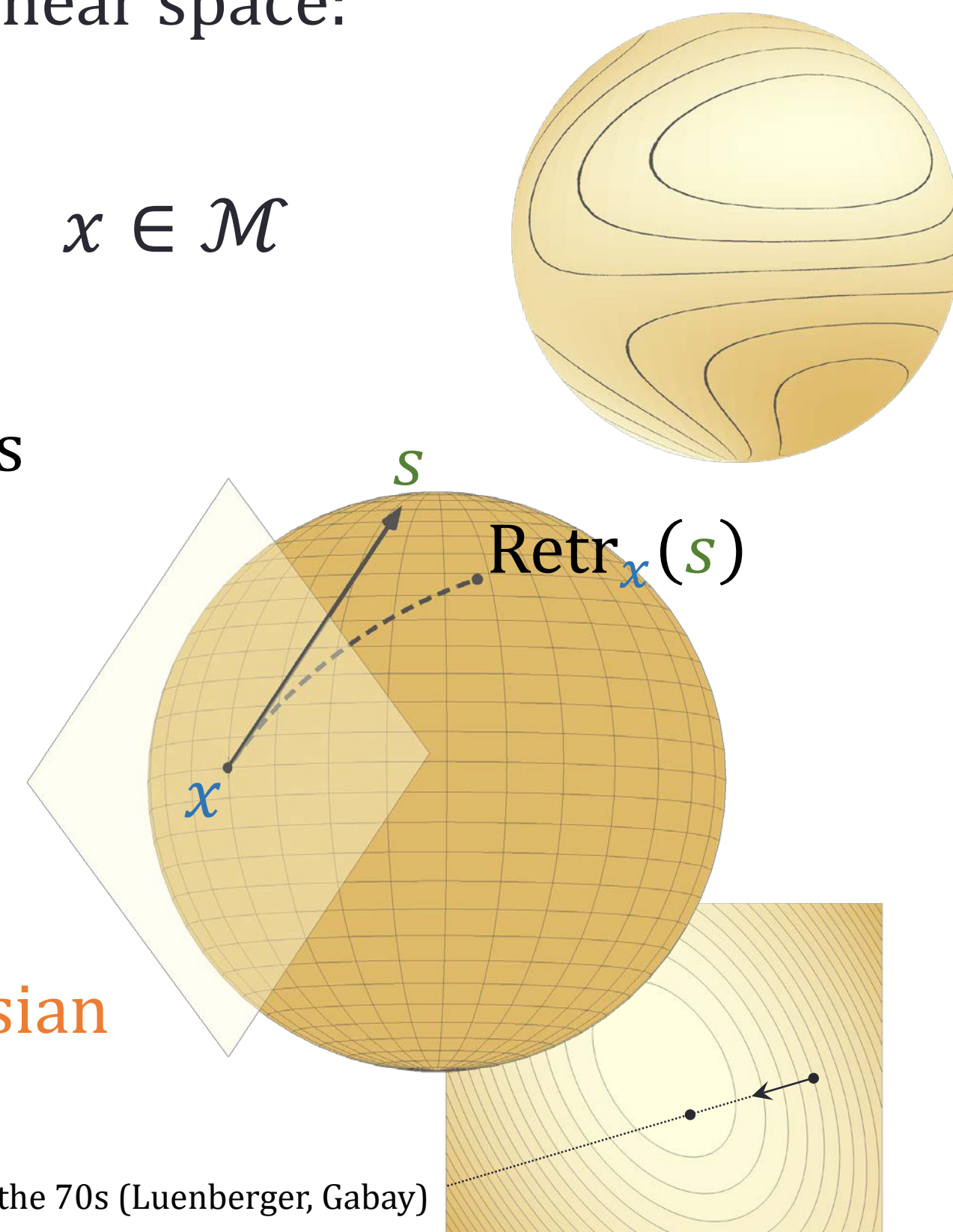
$$\min_x f(x) \quad \text{subject to} \quad x \in \mathcal{M}$$

Tangent spaces: allowed directions
E.g.: $T_x \mathcal{M} = \{s \in \mathbf{R}^n : x^T s = 0\}$

Retractions: tools to move around
E.g.: $\text{Retr}_x(s) = \frac{x+s}{\|x+s\|}$

Riemannian metric: **gradient, Hessian**
E.g.: $\langle s_1, s_2 \rangle_x = s_1^T s_2$

These ideas have been around since the 70s (Luenberger, Gabay)



Algorithms we know and love

$$x_{k+1} = \text{Retr}_{x_k}(s_k)$$

Riemannian gradient descent, trust-regions and cubic regularization:

- $s_k = -\alpha_k \text{grad} f(x_k)$
- $s_k \approx \underset{s \in T_{x_k} \mathcal{M}, \|s\| \leq \Delta_k}{\text{argmin}} f(x_k) + \langle s, \text{grad} f(x_k) \rangle + \frac{1}{2} \langle s, \text{Hess} f(x_k)[s] \rangle$
- $s_k \approx \underset{s \in T_{x_k} \mathcal{M}}{\text{argmin}} f(x_k) + \langle s, \text{grad} f(x_k) \rangle + \frac{1}{2} \langle s, \text{Hess} f(x_k)[s] \rangle + \frac{\sigma_k}{3} \|s\|^3$

Manifolds that matter

Any Cartesian products of all of these:

- Unit norm vectors (spheres)
- Matrices with **orthonormal columns** (Stiefel manifold)
- Subspaces of \mathbf{R}^n of dimension k (Grassmann manifold)
- Fixed-rank matrices (general, symmetric, psd...)
- Low-rank **tensors** (Tucker, tensor train)
- Euclidean distance matrices
- **Rotation** matrices
- Positive probability distributions
- Positive definite matrices
- Many **quotients by group actions**
- ...

Familiar looking bounds

For any $x_0 \in \mathcal{M}$, worst-case iteration complexity:

Riemannian gradient descent

$$O(\varepsilon^{-2}) \text{ for } \|\text{grad} f(x)\| \leq \varepsilon$$

Riemannian trust regions

$$O(\varepsilon^{-2}) \text{ for } \|\text{grad} f(x)\| \leq \varepsilon$$

$$O(\varepsilon^{-3}) \text{ for } \lambda_{\min}(\text{Hess} f(x)) \geq -\varepsilon \text{ too}$$

Riemannian adaptive regularization with cubics

$$O(\varepsilon^{-1.5}) \text{ for small gradient (optimal)}$$

$$O(\varepsilon^{-3}) \text{ for second-order too}$$

RGD and RTR: arXiv:1605.08101, IMA JNA 2019

ARC: arXiv:1806.00065 (see also Zhang & Zhang, arXiv:1805.05565)

Proof for gradient descent

A1 $f(x) \geq f_{\text{low}}$ for all $x \in \mathcal{M}$

A2 $f(\text{Retr}_x(s)) - f(x) - \langle s, \text{grad} f(x) \rangle \leq \frac{L}{2} \|s\|^2$

Algorithm: $x_{k+1} = \text{Retr}_{x_k} \left(-\frac{1}{L} \text{grad} f(x_k) \right)$

Complexity: $\|\text{grad} f(x_k)\| \leq \varepsilon$ with $K \leq 2L(f(x_0) - f_{\text{low}}) \frac{1}{\varepsilon^2}$

$$\text{A2} \Rightarrow f(x_{k+1}) - f(x_k) + \frac{1}{L} \|\text{grad} f(x_k)\|^2 \leq \frac{1}{2L} \|\text{grad} f(x_k)\|^2$$

$$\Rightarrow f(x_k) - f(x_{k+1}) \geq \frac{1}{2L} \|\text{grad} f(x_k)\|^2$$

$$\text{A1} \Rightarrow f(x_0) - f_{\text{low}} \geq \sum_{k=0}^{K-1} f(x_k) - f(x_{k+1}) > \frac{\varepsilon^2}{2L} (K+1)$$

Main assumptions

Mostly the same as in linear spaces:

- f **lower-bounded** on \mathcal{M}
- **Sufficient decrease** per iteration, either in the actual cost function (RGD) or in the model (RTR, ARC)
- **Regularity** assumptions for f ; this is key!
Standard assumptions would be, e.g., Lipschitz gradient.
However, this is uncomfortable on manifolds:

$$\text{dist}(\text{grad} f(x), \text{grad} f(y)) \leq L \text{dist}(x, y)$$

Far easier to compare only scalars:

$$f(\text{Retr}_x(s)) \leq f(x) + \langle s, \text{grad} f(x) \rangle + \frac{L}{2} \|s\|^2$$

Same thing for Lipschitz Hessian: go up one order.

When things are different

To get optimal rates for ARC, we need more work.

Pullback: $\hat{f}_x = f \circ \text{Retr}_x : T_x \mathcal{M} \rightarrow \mathbf{R}$

Solving the subproblem, we make $\text{grad} \hat{f}_x(s)$ small.

For complexity bound, we need $\text{grad} f(\text{Retr}_x(s))$ small.

On linear spaces, they are the same ($\text{Retr}_x(s) = x + s$).

On manifolds,

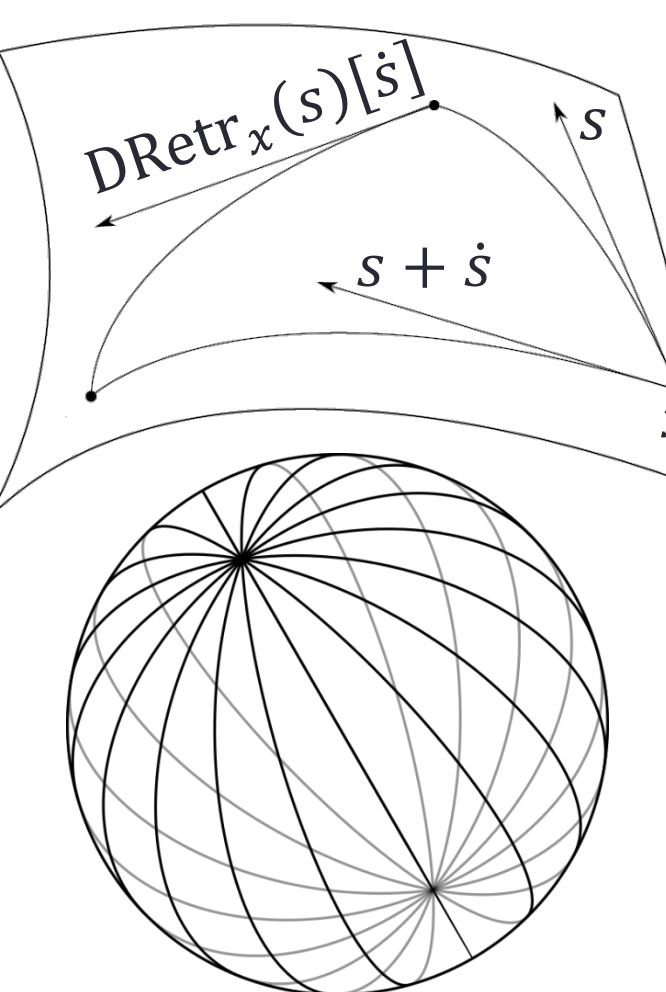
$$\text{grad} \hat{f}_x(s) = (\text{DRetr}_x(s))^{\text{adj}} [\text{grad} f(\text{Retr}_x(s))]$$

Need to control minimum singular value; two tools:

Jacobi field comparison theorem (Lee 1997):

Theorem 11.2. (Jacobi Field Comparison Theorem) Suppose (M, g) is a Riemannian manifold with all sectional curvatures bounded above by a constant C . If γ is a unit speed geodesic in M , and J is any normal Jacobi field along γ such that $J(0) = 0$, then

$$|J(t)| \geq \begin{cases} t |D_t J(0)| & \text{for } 0 \leq t, & \text{if } C = 0; \\ R \sin \frac{t}{R} |D_t J(0)| & \text{for } 0 \leq t \leq \pi R, & \text{if } C = \frac{1}{R^2} > 0; \\ R \sinh \frac{t}{R} |D_t J(0)| & \text{for } 0 \leq t, & \text{if } C = -\frac{1}{R^2} < 0. \end{cases}$$



Maximum theorem (Bergé 1963):

Theorem 2. If ϕ is an upper semi-continuous numerical function in $X \times Y$ and Γ is a u.s.c. mapping of X into Y such that, for each x , $\Gamma x \neq \emptyset$, the numerical function M defined by

$$M(x) = \max \{ \phi(x, y) \mid y \in \Gamma x \}$$

is upper semi-continuous.