

Riemannian Proximal Gradient Methods

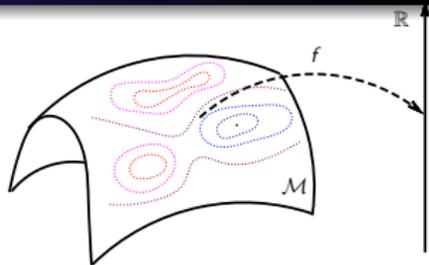
Wen Huang
Xiamen University

Symposium on the Frontiers of Mathematical Optimization Research
Guangxi University July 22, 2020

This is joint work with Ke Wei at Fudan University.

Optimization on Manifolds with Structure:

$$\min_{x \in \mathcal{M}} F(x) = f(x) + g(x),$$

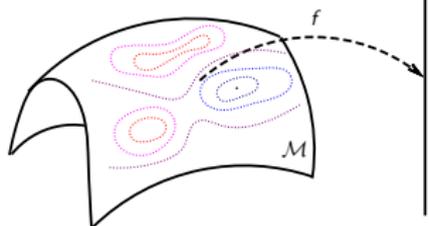


- \mathcal{M} is a Riemannian manifold;
- f is continuously differentiable and may be nonconvex; and
- g is continuous, but may be not differentiable.

¹a penalized version of the ScoTLASS introduced in [JTU03].

Optimization on Manifolds with Structure:

$$\min_{x \in \mathcal{M}} F(x) = f(x) + g(x),$$



- \mathcal{M} is a Riemannian manifold;
- f is continuously differentiable and may be nonconvex; and
- g is continuous, but may be not differentiable.

Applications: sparse PCA, sparse blind deconvolution, sparse low rank image representation, etc [JTU03, GHT15, SQ16, ZLK⁺17]

A sparse PCA optimization model:¹

$$\min_{X \in \text{St}(p, n)} -\text{trace}(X^T A^T A X) + \lambda \|X\|_1,$$

¹a penalized version of the ScoTLASS introduced in [JTU03].

Existing Nonsmooth Optimization on Manifolds

$F : \mathcal{M} \rightarrow \mathbb{R}$ is Lipschitz continuous

- [Huang \(2013\)](#), Gradient sampling method without convergence analysis.
- [Grohs and Hosseini \(2015\)](#), Two ϵ -subgradient-based optimization methods using line search strategy and trust region strategy, respectively. Any limit point is a critical point.
- [Hosseini and Uschmajew \(2017\)](#), Gradient sampling method and any limit point is a critical point.
- [Hosseini and Huang and Yousefpour \(2018\)](#), Merge ϵ -subgradient-based and quasi-Newton ideas and show any limit point is a critical point.

Existing Nonsmooth Optimization on Manifolds

$F : \mathcal{M} \rightarrow \mathbb{R}$ is convex

- [Zhang and Sra \(2016\)](#), subgradient-based method and function value converges to the optimal $O(1/\sqrt{k})$.
- [Ferreira and Oliveira \(2002\)](#) proximal point method, convergence using convexity
[Bento, da Cruz Neto and Oliveira \(2011\)](#), convergence using Kurdyka-Łojasiewicz (KL); and
[Bento, Ferreira and Melo \(2017\)](#), function value converges to the optimal $O(1/k)$ on Hadamard manifold using convexity

Existing Nonsmooth Optimization on Manifolds

$F = f + g$, where f is L-con, and g is non-smooth

- [Chen, Ma, So, and Zhang \(2018\)](#), A proximal gradient method with global convergence
- [Huang and Wei \(2019\)](#), A Riemannian proximal gradient method with convergence rate analyses

A Euclidean Proximal Gradient Method

Optimization with Structure: $\mathcal{M} = \mathbb{R}^n$

$$\min_{x \in \mathbb{R}^n} F(x) = f(x) + g(x), \quad (1)$$

A Euclidean Proximal Gradient Method

Optimization with Structure: $\mathcal{M} = \mathbb{R}^n$

$$\min_{x \in \mathbb{R}^n} F(x) = f(x) + g(x), \quad (1)$$

A proximal gradient method²:

initial iterate: x_0 ,

$$\begin{cases} d_k = \arg \min_{p \in \mathbb{R}^n} \langle \nabla f(x_k), p \rangle + \frac{L}{2} \|p\|_F^2 + g(x_k + p), & \text{(Proximal mapping)} \\ x_{k+1} = x_k + d_k. & \text{(Update iterates)} \end{cases}$$

²The update rule: $x_{k+1} = \arg \min_x \langle \nabla f(x_k), x - x_k \rangle + \frac{L}{2} \|x - x_k\|^2 + g(x)$.

A Euclidean Proximal Gradient Method

Optimization with Structure: $\mathcal{M} = \mathbb{R}^n$

$$\min_{x \in \mathbb{R}^n} F(x) = f(x) + g(x), \quad (1)$$

A proximal gradient method²:

initial iterate: x_0 ,

$$\begin{cases} d_k = \arg \min_{p \in \mathbb{R}^n} \langle \nabla f(x_k), p \rangle + \frac{L}{2} \|p\|_F^2 + g(x_k + p), & \text{(Proximal mapping)} \\ x_{k+1} = x_k + d_k. & \text{(Update iterates)} \end{cases}$$

- $g = 0$: reduce to steepest descent method;

²The update rule: $x_{k+1} = \arg \min_x \langle \nabla f(x_k), x - x_k \rangle + \frac{L}{2} \|x - x_k\|^2 + g(x)$.

A Euclidean Proximal Gradient Method

Optimization with Structure: $\mathcal{M} = \mathbb{R}^n$

$$\min_{x \in \mathbb{R}^n} F(x) = f(x) + g(x), \quad (1)$$

A proximal gradient method²:

initial iterate: x_0 ,

$$\begin{cases} d_k = \arg \min_{p \in \mathbb{R}^n} \langle \nabla f(x_k), p \rangle + \frac{L}{2} \|p\|_F^2 + g(x_k + p), & \text{(Proximal mapping)} \\ x_{k+1} = x_k + d_k. & \text{(Update iterates)} \end{cases}$$

- $g = 0$: reduce to steepest descent method;
- L : greater than the Lipschitz constant of ∇f ;

²The update rule: $x_{k+1} = \arg \min_x \langle \nabla f(x_k), x - x_k \rangle + \frac{L}{2} \|x - x_k\|^2 + g(x)$.

A Euclidean Proximal Gradient Method

Optimization with Structure: $\mathcal{M} = \mathbb{R}^n$

$$\min_{x \in \mathbb{R}^n} F(x) = f(x) + g(x), \quad (1)$$

A proximal gradient method²:

initial iterate: x_0 ,

$$\begin{cases} d_k = \arg \min_{p \in \mathbb{R}^n} \langle \nabla f(x_k), p \rangle + \frac{L}{2} \|p\|_F^2 + g(x_k + p), & \text{(Proximal mapping)} \\ x_{k+1} = x_k + d_k. & \text{(Update iterates)} \end{cases}$$

- $g = 0$: reduce to steepest descent method;
- L : greater than the Lipschitz constant of ∇f ;
- **Proximal mapping: easy to compute;**

²The update rule: $x_{k+1} = \arg \min_x \langle \nabla f(x_k), x - x_k \rangle + \frac{L}{2} \|x - x_k\|^2 + g(x)$.

A Euclidean Proximal Gradient Method

Optimization with Structure: $\mathcal{M} = \mathbb{R}^n$

$$\min_{x \in \mathbb{R}^n} F(x) = f(x) + g(x), \quad (1)$$

A proximal gradient method²:

initial iterate: x_0 ,

$$\begin{cases} d_k = \arg \min_{p \in \mathbb{R}^n} \langle \nabla f(x_k), p \rangle + \frac{L}{2} \|p\|_F^2 + g(x_k + p), & \text{(Proximal mapping)} \\ x_{k+1} = x_k + d_k. & \text{(Update iterates)} \end{cases}$$

- $g = 0$: reduce to steepest descent method;
- L : greater than the Lipschitz constant of ∇f ;
- Proximal mapping: easy to compute;
- **Any limit point is a critical point;**

²The update rule: $x_{k+1} = \arg \min_x \langle \nabla f(x_k), x - x_k \rangle + \frac{L}{2} \|x - x_k\|^2 + g(x)$.

Convergence Rates

Assumption

$\min_{x \in \mathbb{R}^{n \times m}} F(x) = f(x) + g(x)$, with convex f and g ;

- $O(1/k)$ sublinear convergence rate:

$$F(x_k) - F(x_*) \leq C/k, \text{ for a constant } C;$$

Assumption

$\min_{x \in \mathbb{R}^{n \times m}} F(x) = f(x) + g(x)$, with convex f and g ;

- $O(1/k)$ sublinear convergence rate:

$$F(x_k) - F(x_*) \leq C/k, \text{ for a constant } C;$$

- Optimal gradient method: $O(1/k^2)$ [Dar83, Nes83]

Assumption

$\min_{x \in \mathbb{R}^{n \times m}} F(x) = f(x) + g(x)$, with convex f and g ;

- $O(1/k)$ sublinear convergence rate:

$$F(x_k) - F(x_*) \leq C/k, \text{ for a constant } C;$$

- Optimal gradient method: $O(1/k^2)$ [Dar83, Nes83]
- For example: FISTA [BT09]

initial iterate: x_0 and let $y_0 = x_0$, $t_0 = 1$,

$$\begin{cases} d_k = \arg \min_{p \in \mathbb{R}^{n \times m}} \langle \nabla f(y_k), p \rangle + \frac{1}{2} \|p\|_F^2 + g(y_k + p), \\ x_{k+1} = y_k + d_k, \\ t_{k+1} = \frac{1 + \sqrt{4t_k^2 + 1}}{2}, \\ y_{k+1} = x_{k+1} + \frac{t_k - 1}{t_{k+1}} (x_{k+1} - x_k). \end{cases}$$

Assumption

$\min_{x \in \mathbb{R}^n \times m} F(x) = f(x) + g(x)$, with F satisfying the Kurdyka-Łojasiewicz (KL) property with exponent $\theta \in (0, 1]$;

Reference [BST14]:

- Only one accumulation point;
- if $\theta = 1$, then the proximal gradient method terminates in finite steps;
- if $\theta \in [0.5, 1)$, then $\|x_k - x_*\| < C_1 d^k$ for $C_1 > 0$ and $d \in (0, 1)$;
- if $\theta \in (0, 0.5)$, then $\|x_k - x_*\| < C_2 k^{\frac{-1}{1-2\theta}}$ for $C_2 > 0$;

Euclidean proximal mapping

$$d_k = \arg \min_{p \in \mathbb{R}^{n \times m}} \langle \nabla f(x_k), p \rangle + \frac{L}{2} \|p\|_F^2 + g(x_k + p)$$

In the Riemannian setting:

- How to define the proximal mapping?
- Can be solved cheaply?
- Share the same convergence rate?

Euclidean proximal mapping

$$d_k = \arg \min_{p \in \mathbb{R}^{n \times m}} \langle \nabla f(x_k), p \rangle + \frac{L}{2} \|p\|_F^2 + g(x_k + p)$$

A Riemannian proximal mapping [CMMCSZ20]

$$\textcircled{1} \quad \eta_k = \arg \min_{\eta \in \mathbb{T}_{x_k} \mathcal{M}} \langle \nabla f(x_k), \eta \rangle + \frac{L}{2} \|\eta\|_F^2 + g(x_k + \eta);$$

- Only works for embedded submanifold;

²[CMSZ18]: S. Chen, S. Ma, M. C. So, and T. Zhang, Proximal gradient method for nonsmooth optimization over the Stiefel manifold. *SIAM Journal on Optimization*, 30(1):210-239, 2020

Euclidean proximal mapping

$$d_k = \arg \min_{p \in \mathbb{R}^{n \times m}} \langle \nabla f(x_k), p \rangle + \frac{L}{2} \|p\|_F^2 + g(x_k + p)$$

A Riemannian proximal mapping [CMMCSZ20]

$$\textcircled{1} \quad \eta_k = \arg \min_{\eta \in T_{x_k} \mathcal{M}} \langle \nabla f(x_k), \eta \rangle + \frac{L}{2} \|\eta\|_F^2 + g(x_k + \eta);$$

- Only works for embedded submanifold;
- Proximal mapping is defined in tangent space;

²[CMSZ18]: S. Chen, S. Ma, M. C. So, and T. Zhang, Proximal gradient method for nonsmooth optimization over the Stiefel manifold. *SIAM Journal on Optimization*, 30(1):210-239, 2020

Euclidean proximal mapping

$$d_k = \arg \min_{p \in \mathbb{R}^{n \times m}} \langle \nabla f(x_k), p \rangle + \frac{L}{2} \|p\|_F^2 + g(x_k + p)$$

A Riemannian proximal mapping [CMMCSZ20]

$$1 \quad \eta_k = \arg \min_{\eta \in T_{x_k} \mathcal{M}} \langle \nabla f(x_k), \eta \rangle + \frac{L}{2} \|\eta\|_F^2 + g(x_k + \eta);$$

- Only works for embedded submanifold;
- Proximal mapping is defined in tangent space;
- **Convex programming;**

²[CMSZ18]: S. Chen, S. Ma, M. C. So, and T. Zhang, Proximal gradient method for nonsmooth optimization over the Stiefel manifold. *SIAM Journal on Optimization*, 30(1):210-239, 2020

Euclidean proximal mapping

$$d_k = \arg \min_{p \in \mathbb{R}^{n \times m}} \langle \nabla f(x_k), p \rangle + \frac{L}{2} \|p\|_F^2 + g(x_k + p)$$

ManPG [CMMCSZ20]

$$\bullet \eta_k = \arg \min_{\eta \in T_{x_k} \mathcal{M}} \langle \nabla f(x_k), \eta \rangle + \frac{L}{2} \|\eta\|_F^2 + g(x_k + \eta);$$

- Only works for embedded submanifold;
- Proximal mapping is defined in tangent space;
- Convex programming;
- Solved efficiently for the Stiefel manifold by a semi-Newton algorithm [XLWZ18];

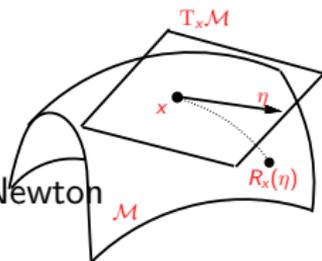
Euclidean proximal mapping

$$d_k = \arg \min_{p \in \mathbb{R}^{n \times m}} \langle \nabla f(x_k), p \rangle + \frac{L}{2} \|p\|_F^2 + g(x_k + p)$$

ManPG [CMMCSZ20]

- 1 $\eta_k = \arg \min_{\eta \in T_{x_k} \mathcal{M}} \langle \nabla f(x_k), \eta \rangle + \frac{L}{2} \|\eta\|_F^2 + g(x_k + \eta);$
- 2 $x_{k+1} = R_{x_k}(\alpha_k \eta_k)$ with an appropriate step size α_k ;

- Only works for embedded submanifold;
- Proximal mapping is defined in tangent space;
- Convex programming;
- Solved efficiently for the Stiefel manifold by a semi-Newton algorithm [XLWZ18];
- Step size 1 is not necessary decreasing;



Euclidean proximal mapping

$$d_k = \arg \min_{p \in \mathbb{R}^{n \times m}} \langle \nabla f(x_k), p \rangle + \frac{L}{2} \|p\|_F^2 + g(x_k + p)$$

ManPG [CMMCSZ20]

- 1 $\eta_k = \arg \min_{\eta \in T_{x_k}} \mathcal{M} \langle \nabla f(x_k), \eta \rangle + \frac{L}{2} \|\eta\|_F^2 + g(x_k + \eta);$
 - 2 $x_{k+1} = R_{x_k}(\alpha_k \eta_k)$ with an appropriate step size α_k ;
- Convergence to a stationary point;

Euclidean proximal mapping

$$d_k = \arg \min_{p \in \mathbb{R}^{n \times m}} \langle \nabla f(x_k), p \rangle + \frac{L}{2} \|p\|_F^2 + g(x_k + p)$$

ManPG [CMMCSZ20]

- 1 $\eta_k = \arg \min_{\eta \in T_{x_k}} \mathcal{M} \langle \nabla f(x_k), \eta \rangle + \frac{L}{2} \|\eta\|_F^2 + g(x_k + \eta);$
 - 2 $x_{k+1} = R_{x_k}(\alpha_k \eta_k)$ with an appropriate step size α_k ;
- Convergence to a stationary point;
 - **No convergence rate analysis;**

New Riemannian Proximal Gradient Methods

GOAL: Develop a Riemannian proximal gradient method with convergence rate analysis and good numerical performance for some instances

New Riemannian Proximal Gradient Methods

GOAL: Develop a Riemannian proximal gradient method with convergence rate analysis and good numerical performance for some instances

A New Riemannian Proximal Gradient Method

- 1 $\eta_k = \arg \min_{\eta \in T_{x_k}} \mathcal{M} \underbrace{\langle \nabla f(x_k), \eta \rangle_{x_k} + \frac{L}{2} \|\eta\|_{x_k}^2}_{\text{Riemannian metric}} + g(\underbrace{R_{x_k}(\eta)}_{\text{replace } x_k + \eta});$
- 2 $x_{k+1} = R_{x_k}(\eta_k);$

- General framework for Riemannian optimization;

New Riemannian Proximal Gradient Methods

GOAL: Develop a Riemannian proximal gradient method with convergence rate analysis and good numerical performance for some instances

A New Riemannian Proximal Gradient Method

$$\textcircled{1} \quad \eta_k = \arg \min_{\eta \in T_{x_k}} \underbrace{\mathcal{M} \langle \nabla f(x_k), \eta \rangle_{x_k} + \frac{L}{2} \|\eta\|_{x_k}^2}_{\text{Riemannian metric}} + g(\underbrace{R_{x_k}(\eta)}_{\text{replace } x_k + \eta});$$
$$\textcircled{2} \quad x_{k+1} = R_{x_k}(\eta_k);$$

- General framework for Riemannian optimization;
- Step size can be fixed to be 1;

Assumptions and Convergence Result

Assumption:

- 1 The function F is bounded from below and the sublevel set $\Omega_{x_0} = \{x \in \mathcal{M} \mid F(x) \leq F(x_0)\}$ is compact;

This assumption hold if, for example, F is continuous and \mathcal{M} is compact.

$$\min_{X \in \text{St}(p,n)} -\text{trace}(X^T A^T A X) + \lambda \|X\|_1,$$

Assumptions and Convergence Result

Assumption:

- 1 The function F is bounded from below and the sublevel set $\Omega_{x_0} = \{x \in \mathcal{M} \mid F(x) \leq F(x_0)\}$ is compact;
 - 2 The function f is L -retraction-smooth with respect to the retraction R in the sublevel set Ω_{x_0} .
-

Definition

A function $h : \mathcal{M} \rightarrow \mathbb{R}$ is called L -retraction-smooth with respect to a retraction R in $\mathcal{N} \subseteq \mathcal{M}$ if for any $x \in \mathcal{N}$ and any $\mathcal{S}_x \subseteq T_x \mathcal{M}$ such that $R_x(\mathcal{S}_x) \subseteq \mathcal{N}$, we have that

$$h(R_x(\eta)) \leq h(x) + \langle \text{grad } h(x), \eta \rangle_x + \frac{L}{2} \|\eta\|_x^2, \quad \forall \eta \in \mathcal{S}_x.$$

Assumptions and Convergence Result

Assumption:

- 1 The function F is bounded from below and the sublevel set $\Omega_{x_0} = \{x \in \mathcal{M} \mid F(x) \leq F(x_0)\}$ is compact;
- 2 The function f is L -retraction-smooth with respect to the retraction R in the sublevel set Ω_{x_0} .

if the following conditions hold, then f is L -retraction-smooth with respect to the retraction R in the manifold \mathcal{M} [BAC18, Lemma 2.7]

- \mathcal{M} is a compact Riemannian submanifold of a Euclidean space \mathbb{R}^n ;
- the retraction R is globally defined;
- $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is L -smooth in the convex hull of \mathcal{M} ;

$$\min_{X \in \text{St}(p, n)} -\text{trace}(X^T A^T A X) + \lambda \|X\|_1,$$

Assumptions and Convergence Result

Assumption:

- 1 The function F is bounded from below and the sublevel set $\Omega_{x_0} = \{x \in \mathcal{M} \mid F(x) \leq F(x_0)\}$ is compact;
 - 2 The function f is L -retraction-smooth with respect to the retraction R in the sublevel set Ω_{x_0} .
-

Theoretical results:

- For any accumulation point x_* of $\{x_k\}$, x_* is a stationary point, i.e., $0 \in \partial F(x_*)$.

Assumptions and Convergence Rate

Additional Assumptions:

- f and g are retraction-convex in $\Omega \supseteq \Omega_{x_0}$;
-

Definition

A function $h : \mathcal{M} \rightarrow \mathbb{R}$ is called retraction-convex with respect to a retraction R in $\mathcal{N} \subseteq \mathcal{M}$ if for any $x \in \mathcal{N}$ and any $\mathcal{S}_x \subseteq T_x \mathcal{M}$ such that $R_x(\mathcal{S}_x) \subseteq \mathcal{N}$, there exists a tangent vector $\zeta \in T_x \mathcal{M}$ such that $q_x = h \circ R_x$ satisfies

$$q_x(\eta) \geq q_x(\xi) + \langle \zeta, \eta - \xi \rangle_x \quad \forall \eta, \xi \in \mathcal{S}_x. \quad (2)$$

Note that $\zeta = \text{grad } q_x(\xi)$ if h is differentiable; otherwise, ζ is any subgradient of q_x at ξ .

Assumptions and Convergence Rate

Additional Assumptions:

- f and g are retraction-convex in $\Omega \supseteq \Omega_{x_0}$;
-

Lemma

Given $x \in \mathcal{M}$ and a twice continuously differentiable function $h : \mathcal{M} \rightarrow \mathbb{R}$, if one of the following conditions holds:

- Hess h is positive definite at x , and the retraction is second order;
- The manifold \mathcal{M} is an embedded submanifold of \mathbb{R}^n endowed with the Euclidean metric; \mathcal{W} is an open subset of \mathbb{R}^n ; $x \in \mathcal{W}$;
 $h : \mathcal{W} \subset \mathbb{R}^n \rightarrow \mathbb{R}$ is a μ -strongly convex function in the Euclidean setting for a sufficient large μ ; the retraction is second order;

then there exists a neighborhood of x , denoted by \mathcal{N}_x , such that the function $h : \mathcal{M} \rightarrow \mathbb{R}$ is retraction-convex in \mathcal{N}_x .

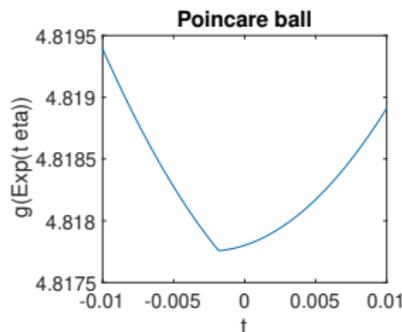
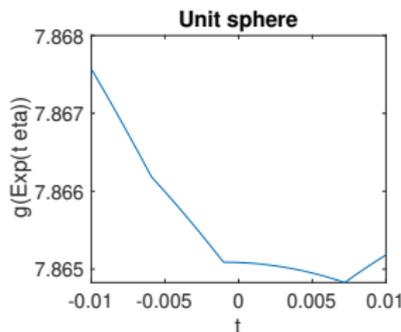
Assumptions and Convergence Rate

Additional Assumptions:

- f and g are retraction-convex in $\Omega \supseteq \Omega_{x_0}$;

Nonsmooth? Example: $g(x) = \|x\|_1$ with exponential mapping

- unit sphere: $\{x \in \mathbb{R}^n \mid x^T x = 1\}$, $n = 100$
- Poincaré ball model [GBH18]: $\{x \in \mathbb{R}^n \mid x^T x < 1\}$, $n = 100$
- $g(\text{Exp}_x(t\eta_x))$ versus t



[GBH18] Ganea et al., Hyperbolic entailment cones for learning hierarchical embedding, ICML, 2018.

Assumptions and Convergence Rate

Additional Assumptions:

- f and g are retraction-convex in $\Omega \supseteq \Omega_{x_0}$;
- Retraction approximately satisfies the triangle relation in Ω : for all $x, y, z \in \Omega$,

$$|\|\xi_x - \eta_x\|_x^2 - \|\zeta_y\|_y^2| \leq \kappa \|\eta_x\|_x^2, \text{ for a constant } \kappa$$

where $\eta_x = R_x^{-1}(y)$, $\xi_x = R_x^{-1}(z)$, $\zeta_y = R_y^{-1}(z)$.

-
- In the Euclidean setting: $\eta_x = R_x^{-1}(y) = y - x$, $\xi_x = R_x^{-1}(z) = z - x$, $\zeta_y = R_y^{-1}(z) = z - y$:

$$\xi_x - \eta_x = (z - x) - (y - x) = z - y = \zeta_y.$$

- Holds on the unit sphere.

Assumptions and Convergence Rate

Additional Assumptions:

- f and g are retraction-convex in $\Omega \supseteq \Omega_{x_0}$;
- Retraction approximately satisfies the triangle relation in Ω : for all $x, y, z \in \Omega$,

$$|\|\xi_x - \eta_x\|_x^2 - \|\zeta_y\|_y^2| \leq \kappa \|\eta_x\|_x^2, \text{ for a constant } \kappa$$

where $\eta_x = R_x^{-1}(y)$, $\xi_x = R_x^{-1}(z)$, $\zeta_y = R_y^{-1}(z)$.

Table: Exponential mapping on the Stiefel manifold with the Euclidean metric $\langle \eta_x, \xi_x \rangle_x = \text{trace}(\eta_x^T \xi_x)$. Left = $|\|\xi_x - \eta_x\|_x^2 - \|\zeta_y\|_y^2|$

$(n, p) = (10, 1)$		$(n, p) = (10, 4)$		$(n, p) = (10, 10)$	
$\ \eta_x\ $	Left	$\ \eta_x\ $	Left	$\ \eta_x\ $	Left
5.00 ₋₂	7.83 ₋₅	5.00 ₋₂	1.83 ₋₅	5.00 ₋₂	2.14 ₋₆
2.50 ₋₂	1.80 ₋₅	2.50 ₋₂	4.27 ₋₆	2.50 ₋₂	4.72 ₋₇
1.25 ₋₂	4.25 ₋₆	1.25 ₋₂	1.01 ₋₆	1.25 ₋₂	1.11 ₋₇
6.25 ₋₃	1.03 ₋₆	6.25 ₋₃	2.46 ₋₇	6.25 ₋₃	2.68 ₋₈

Assumptions and Convergence Rate

Additional Assumptions:

- f and g are retraction-convex in $\Omega \supseteq \Omega_{x_0}$;
- Retraction approximately satisfies the triangle relation in Ω : for all $x, y, z \in \Omega$,

$$|\|\xi_x - \eta_x\|_x^2 - \|\zeta_y\|_y^2| \leq \kappa \|\eta_x\|_x^2, \text{ for a constant } \kappa$$

where $\eta_x = R_x^{-1}(y)$, $\xi_x = R_x^{-1}(z)$, $\zeta_y = R_y^{-1}(z)$.

Table: Exponential mapping on the Stiefel manifold with the canonical metric $\langle \eta_x, \xi_x \rangle_x = \text{trace}(\eta_x^T (I - XX^T/2) \xi_x)$. Left = $|\|\xi_x - \eta_x\|_x^2 - \|\zeta_y\|_y^2|$

$(n, p) = (10, 2)$		$(n, p) = (10, 4)$		$(n, p) = (10, 9)$	
$\ \eta_x\ $	Left	$\ \eta_x\ $	Left	$\ \eta_x\ $	Left
5.00 ₋₂	3.55 ₋₅	5.00 ₋₂	1.15 ₋₅	5.00 ₋₂	8.39 ₋₆
2.50 ₋₂	8.06 ₋₆	2.50 ₋₂	2.58 ₋₆	2.50 ₋₂	1.89 ₋₆
1.25 ₋₂	1.90 ₋₆	1.25 ₋₂	6.08 ₋₇	1.25 ₋₂	4.45 ₋₇
6.25 ₋₃	4.61 ₋₇	6.25 ₋₃	1.47 ₋₇	6.25 ₋₃	1.08 ₋₇

Assumptions and Convergence Rate

Additional Assumptions:

- f and g are retraction-convex in $\Omega \supseteq \Omega_{x_0}$;
- Retraction approximately satisfies the triangle relation in Ω : for all $x, y, z \in \Omega$,

$$\left| \|\xi_x - \eta_x\|_x^2 - \|\zeta_y\|_y^2 \right| \leq \kappa \|\eta_x\|_x^2, \text{ for a constant } \kappa$$

where $\eta_x = R_x^{-1}(y)$, $\xi_x = R_x^{-1}(z)$, $\zeta_y = R_y^{-1}(z)$.

Theoretical results:

- Convergence rate $O(1/k)$:

$$F(x_k) - F(x_*) \leq \frac{1}{k} \left(\frac{L}{2} \|R_{x_0}^{-1}(x_*)\|_{x_0}^2 + \frac{L\kappa C}{2} (F(x_0) - F(x_*)) \right).$$

Riemannian FISTA Method with $O(1/k^2)$?

FISTA initial iterate: x_0 and let $y_0 = x_0$, $t_0 = 1$

$$\textcircled{1} \quad d_k = \arg \min_{p \in \mathbb{R}^{n \times m}} \langle \nabla f(y_k), p \rangle + \frac{L}{2} \|p\|_F^2 + g(y_k + p)$$

$$\textcircled{2} \quad x_{k+1} = y_k + d_k$$

$$\textcircled{3} \quad t_{k+1} = \frac{1 + \sqrt{4t_k^2 + 1}}{2}$$

$$\textcircled{4} \quad y_{k+1} = x_{k+1} + \frac{t_k - 1}{t_{k+1}} (x_{k+1} - x_k)$$

Possible Riemannian generalizations:

- Step 1: Riemannian proximal mapping
- Step 2: Retraction
- Step 4: multiple generalizations

Difficulties for $O(1/k^2)$ convergence rate, e.g.,

$$\left| \|\omega_x + \xi_x - \eta_x\|_x^2 - \|\omega_x + \zeta_y\|_y^2 \right| \leq \kappa \|\eta_x\|_x^2, \text{ for a constant } \kappa$$

Assumptions and Local Convergence Result

Assumption:

- 1 Assumptions for the global convergence

-
- 1 The function F is bounded from below and the sublevel set $\Omega_{x_0} = \{x \in \mathcal{M} \mid F(x) \leq F(x_0)\}$ is compact;
 - 2 The function f is L -retraction-smooth with respect to the retraction R in the sublevel set Ω_{x_0} .

$$\min_{X \in \text{St}(p,n)} -\text{trace}(X^T A^T A X) + \lambda \|X\|_1,$$

Assumptions and Local Convergence Result

Assumption:

- 1 Assumptions for the global convergence
- 2 f is locally Lipschitz continuously differentiable

Definition ([AMS08, 7.4.3])

A function f on \mathcal{M} is Lipschitz continuously differentiable if it is differentiable and if there exists β_1 such that, for all x, y in \mathcal{M} with $\text{dist}(x, y) < i(\mathcal{M})$, it holds that

$$\|\mathcal{P}_\gamma^{0 \leftarrow 1} \text{grad } f(y) - \text{grad } f(x)\|_x \leq \beta_1 \text{dist}(x, y),$$

where γ is the unique minimizing geodesic with $\gamma(0) = x$ and $\gamma(1) = y$.

Assumptions and Local Convergence Result

Assumption:

- 1 Assumptions for the global convergence
- 2 f is locally Lipschitz continuously differentiable

If f is smooth and the manifold \mathcal{M} is compact, then the function f is Lipschitz continuously differentiable. [AMS08, Proposition 7.4.5 and Corollary 7.4.6].

$$\min_{X \in \text{St}(p,n)} -\text{trace}(X^T A^T A X) + \lambda \|X\|_1,$$

Assumptions and Local Convergence Result

Assumption:

- 1 Assumptions for the global convergence
 - 2 f is locally Lipschitz continuously differentiable
 - 3 F satisfies the Riemannian KL property [BCNO11]
-

Definition

A continuous function $f : \mathcal{M} \rightarrow \mathbb{R}$ is said to have the Riemannian KL property at $x \in \mathcal{M}$ if and only if there exists $\varepsilon \in (0, \infty]$, a neighborhood $U \subset \mathcal{M}$ of x , and a continuous concave function $\varsigma : [0, \varepsilon] \rightarrow [0, \infty)$ such that

- $\varsigma(0) = 0$, ς is C^1 on $(0, \varepsilon)$, and $\varsigma' > 0$ on $(0, \eta)$,
- For every $y \in U$ with $f(x) < f(y) < f(x) + \varepsilon$, we have

$$\varsigma'(f(y) - f(x)) \operatorname{dist}(0, \partial f(y)) \geq 1,$$

where $\operatorname{dist}(0, \partial f(y)) = \inf\{\|v\|_y : v \in \partial f(y)\}$ and ∂ denotes the Riemannian generalized subdifferential. The function ς is called the desingularising function.

Assumptions and Local Convergence Result

Assumption:

- 1 Assumptions for the global convergence
 - 2 f is locally Lipschitz continuously differentiable
 - 3 F satisfies the Riemannian KL property [BCNO11]
-

Theoretical results:

- it holds that

$$\sum_{k=0}^{\infty} \text{dist}(x_k, x_{k+1}) < \infty.$$

Therefore, there exists only a unique accumulation point.

Assumptions and Local Convergence Result

Assumption:

- 1 Assumptions for the global convergence
 - 2 f is locally Lipschitz continuously differentiable
 - 3 F satisfies the Riemannian KL property [BCNO11]
-

Theoretical results:

- If the desingularising function has the form $\varsigma(t) = \frac{C}{\theta} t^\theta$ for $C > 0$ and $\theta \in (0, 1]$ for all $x \in \Omega_{x_0}$, then
 - if $\theta = 1$, then the Riemannian proximal gradient method terminates in finite steps;
 - if $\theta \in [0.5, 1)$, then $\|x_k - x_*\| < C_1 d^k$ for $C_1 > 0$ and $d \in (0, 1)$;
 - if $\theta \in (0, 0.5)$, then $\|x_k - x_*\| < C_2 k^{\frac{-1}{1-2\theta}}$ for $C_2 > 0$;

How to verify if a function satisfies the Riemannian KL property?

Theorem

Given $x \in \mathcal{M}$, let (ϕ, \mathcal{U}) denote a chart of \mathcal{M} covering x , i.e., $x \in \mathcal{U}$. We assume that $F \circ \phi^{-1} : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies the Euclidean KL property at $\phi(x)$ with the desingularising function $\tilde{\zeta}_x$, then F satisfies the Riemannian KL property at x with the desingularising function $\tilde{\zeta}_x/C_x$, where C_x is a constant.

Similar result is given in [BCNO11]: F is a \mathcal{C} -function $\implies F$ satisfies the Riemannian KL property.

Riemannian KL property

Semialgebraic sets, mappings, and functions

Definition (Semialgebraic sets, mappings and functions)

- 1 A subset \mathcal{S} of \mathbb{R}^n is called semialgebraic if there exists a finite number of polynomial function $g_{ij}, h_{ij} : \mathbb{R}^n \rightarrow \mathbb{R}$ such that

$$\mathcal{S} = \cup_{j=1}^p \cap_{i=1}^q \{u \in \mathbb{R}^n \mid g_{ij}(u) = 0 \text{ and } h_{ij}(u) < 0\}.$$

- 2 Let $\mathcal{A} \subseteq \mathbb{R}^m$ and $\mathcal{B} \subseteq \mathbb{R}^n$ be two semialgebraic sets. A mapping $\gamma : \mathcal{A} \rightarrow \mathcal{B}$ is semialgebraic if its graph is semialgebraic in \mathbb{R}^{m+n} . If $n = 1$, then the mapping is also called a semialgebraic function.

Continuous semialgebraic functions satisfy the Euclidean KL property with desingularising function in the form of $\zeta(t) = \frac{C}{\theta} t^\theta$, where $\theta \in (0, 1]$ and $C > 0$.

Riemannian KL property on the Stiefel manifold

Restriction of a semialgebraic Function onto Stiefel manifold satisfies the Riemannian KL property

- 1 For any point $x \in \text{St}(p, n)$, construct a chart (ϕ, \mathcal{U}) such that $x \in \mathcal{U}$ and ϕ is a semialgebraic mapping
- 2 Inverse of ϕ is semialgebraic mapping
- 3 The composition function $f \circ \phi^{-1}$ is a semialgebraic function
- 4 $f \circ \phi^{-1}$ satisfies the Euclidean KL property with desingularising function $\varsigma(t) = \frac{c}{\theta} t^\theta$
- 5 $f : \mathcal{M} \rightarrow \mathbb{R}$ satisfies the Riemannian KL property with desingularising function $\varsigma(t) = \frac{\tilde{c}}{\theta} t^\theta$

Restriction of a semialgebraic Function onto Stiefel manifold satisfies the Riemannian KL property

- 1 For any point $x \in \text{St}(p, n)$, construct a chart (ϕ, \mathcal{U}) such that $x \in \mathcal{U}$ and ϕ is a semialgebraic mapping
- 2 Inverse of ϕ is semialgebraic mapping
- 3 The composition function $f \circ \phi^{-1}$ is a semialgebraic function
- 4 $f \circ \phi^{-1}$ satisfies the Euclidean KL property with desingularising function $\varsigma(t) = \frac{c}{\theta} t^\theta$
- 5 $f : \mathcal{M} \rightarrow \mathbb{R}$ satisfies the Riemannian KL property with desingularising function $\varsigma(t) = \frac{\tilde{c}}{\theta} t^\theta$

$$\min_{X \in \text{St}(p, n)} -\text{trace}(X^T A^T A X) + \lambda \|X\|_1,$$

$$\min_{X \in \text{St}(p,n)} -\text{trace}(X^T A^T A X) + \lambda \|X\|_1,$$

Therefore, all assumptions for global convergence and local convergence hold for the two sparse PCA models.

- All accumulation points are critical
- Local convergence:
 - Accumulation point is unique;
 - if $\theta = 1$, then the method terminates in finite steps;
 - if $\theta \in [0.5, 1)$, then $\|x_k - x_*\| < C_1 d^k$ for $C_1 > 0$ and $d \in (0, 1)$;
 - if $\theta \in (0, 0.5)$, then $\|x_k - x_*\| < C_2 k^{\frac{-1}{1-2\theta}}$ for $C_2 > 0$;

$$\eta_x = \arg \min_{\eta \in T_x \mathcal{M}} \ell_x(\eta) := \langle \nabla f(x), \eta \rangle_x + \frac{L}{2} \|\eta\|_x^2 + g(R_x(\eta))$$

In some cases, the subproblem can be solved by exploiting the structure of the manifold;

Solving the Riemannian Proximal Mapping

initial iterate: $\eta_0 \in T_x \mathcal{M}$, $\sigma \in (0, 1)$, $k = 0$;

① $y_k = R_x(\eta_k)$;

② Compute

$$\xi_k^* = \arg \min_{\xi \in T_{y_k} \mathcal{M}} \langle \mathcal{T}_{R_{\eta_k}}^{-\sharp}(\text{grad } f(x) + \tilde{L}\eta_k), \xi \rangle_x + \frac{\tilde{L}}{4} \|\xi\|_F^2 + g(y_k + \xi);$$

③ Find $\alpha > 0$ such that $\ell_x(\eta_k + \alpha \mathcal{T}_{R_{\eta_k}}^{-1} \xi_k^*) < \ell_x(\eta_k) - \sigma \alpha \|\xi_k^*\|_x^2$;

④ $\eta_{k+1} = \eta_k + \alpha \mathcal{T}_{R_{\eta_k}}^{-1} \xi_k^*$, $k \leftarrow k + 1$ and goto Step 1;

Above algorithm is used if the ambient space is \mathbb{R}^n

An application of [CMMCSZ20] if $R_x^{-1}(y)$ exists.

Two sparse PCA models:

- first model: [GHT15]

$$\min_{X \in OB(p,n)} \|X^T A^T A X - D^2\|_F^2 + \lambda \|X\|_1,$$

where $A \in \mathbb{R}^{m \times n}$ is a data matrix, D is the diagonal matrix with dominant singular values of A ,

$$OB(p,n) = \{X \in \mathbb{R}^{n \times p} \mid \text{diag}(X^T X) = I_p\}, \quad p \leq m;$$

- second model

$$\min_{X \in \text{St}(p,n)} -\text{trace}(X^T A^T A X) + \lambda \|X\|_1.$$

Numerical Experiments

Table: An average result of 10 random tests. $n = 128$, $m = 20$, $r = 4$.
 $\delta = (L\|x_{k+1} - x_k\|)^2$. The subscript k indicates a scale of 10^k .

λ	Algo	iter	time	f	δ	spar.	navar
3	ManPG	11791	1.40	8.33_1	5.11_{-6}	0.54	0.86
	RPG	11679	0.94	8.33_1	5.11_{-6}	0.54	0.86
	ManPG-Ada	1398	0.30	8.33_1	1.67_{-3}	0.54	0.86
	A-ManPG	273	0.09	8.33_1	9.19_{-4}	0.54	0.86
	A-RPG	263	0.06	8.33_1	1.12_{-3}	0.54	0.86

- **ManPG**: the method in [CMMCSZ20];
- **RPG**: the new Riemannian proximal gradient without acceleration;
- **A-ManPG**: Use similar technique to accelerate ManPG;
- **A-RPG**: the new Riemannian proximal gradient with acceleration;

Numerical Experiments

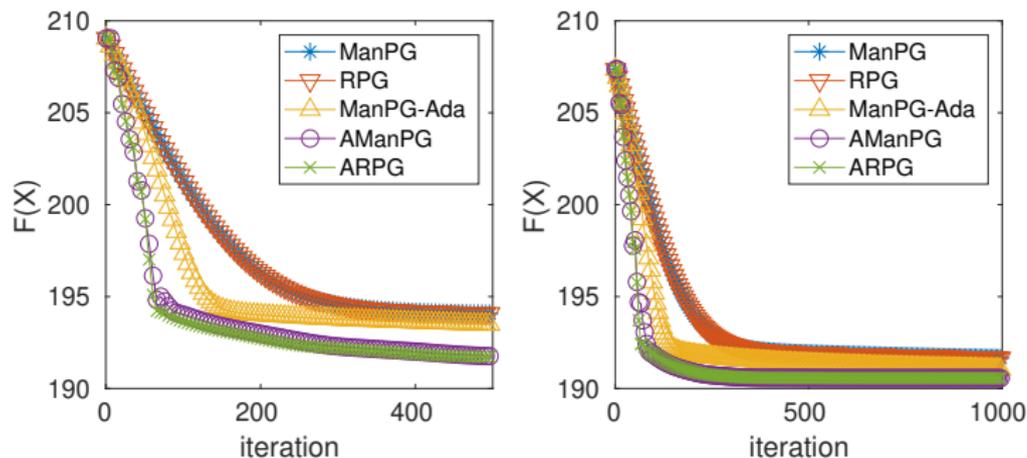


Figure: Two typical runs of ManPG, RPG, A-ManPG, and A-RPG for the Sparse PCA problem. $n = 1024$, $p = 4$, $\lambda = 2$, $m = 20$.

Sparse PCA problem

$$\min_{X \in \text{St}(p,n)} -\text{trace}(X^T A^T A X) + \lambda \|X\|_1,$$

where $A \in \mathbb{R}^{m \times n}$ is a data matrix.

Numerical Experiments

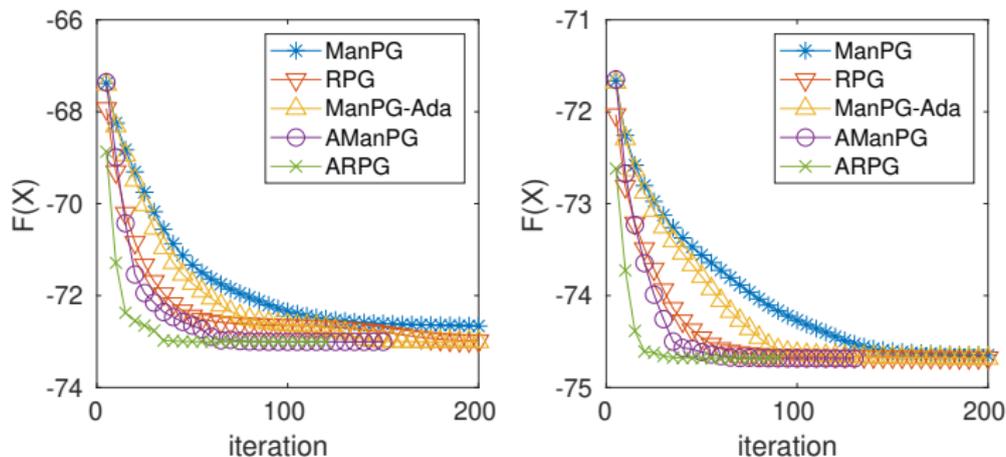


Figure: Two typical runs of ManPG, RPG, A-ManPG, and A-RPG for the Sparse PCA problem. $n = 1024$, $p = 4$, $\lambda = 2$, $m = 20$.

- Propose a Riemannian proximal gradient method;
- Global convergence to critical points
- $O(1/k)$ convergence rate using retraction-convexity
- Local convergence rate using Riemannian KL property
- Retraction of a semialgebraic function onto the Stiefel manifold satisfies the Riemannian KL property
- Apply the methods to sparse PCA problems on the oblique manifold and the Stiefel manifold;

- Wen Huang, Ke Wei, Riemannian Proximal Gradient Methods, arxiv:1909.06065, 2019

Thank you

References I



P.-A. Absil, R. Mahony, and R. Sepulchre.
Optimization algorithms on matrix manifolds.
Princeton University Press, Princeton, NJ, 2008.



Nicolas Boumal, P-A Absil, and Coralia Cartis.
Global rates of convergence for nonconvex optimization on manifolds.
IMA Journal of Numerical Analysis, 39(1):1–33, 02 2018.



G. C. Bento, J. X. Cruz Neto, and P. R. Oliveira.
Convergence of inexact descent methods for nonconvex optimization on Riemannian manifold.
arXiv preprint arXiv:1103.4828, 2011.



Jérôme Bolte, Shoham Sabach, and Marc Teboulle.
Proximal alternating linearized minimization for nonconvex and nonsmooth problems.
Mathematical Programming, 146(1-2):459–494, 2014.



A. Beck and M. Teboulle.
A fast iterative shrinkage-thresholding algorithm for linear inverse problems.
SIAM Journal on Imaging Sciences, 2(1):183–202, January 2009.
doi:10.1137/080716542.



Shixiang Chen, Shiqian Ma, Anthony Man-Cho So, and Tong Zhang.
Proximal gradient method for nonsmooth optimization over the Stiefel manifold.
SIAM Journal on Optimization, 30(1):210–239, 2020.



John Darzentas.
Problem Complexity and Method Efficiency in Optimization.
1983.



Octavian Eugen Ganea, Gary Becigneul, and Thomas Hofmann.
Hyperbolic entailment cones for learning hierarchical embeddings.
35th International Conference on Machine Learning, ICML 2018, 4:2661–2673, 2018.



Matthieu Genicot, Wen Huang, and Nickolay T. Trendafilov.
Weakly correlated sparse components with nearly orthonormal loadings.
In *Geometric Science of Information*, pages 484–490, 2015.



Ian T. Jolliffe, Nickolay T. Trendafilov, and Mudassir Uddin.
A modified principal component technique based on the Lasso.
Journal of Computational and Graphical Statistics, 12(3):531–547, 2003.



Y. E. Nesterov.
A method for solving the convex programming problem with convergence rate $O(1/k^2)$.
Dokl. Akad. Nauk SSSR (In Russian), 269:543–547, 1983.



J. Shi and C. Qi.
Low-rank sparse representation for single image super-resolution via self-similarity learning.
In *2016 IEEE International Conference on Image Processing (ICIP)*, pages 1424–1428, Sep. 2016.



Xiantao Xiao, Yongfeng Li, Zaiwen Wen, and Liwei Zhang.
A regularized semi-smooth newton method with projection steps for composite convex programs.
Journal of Scientific Computing, 76(1):364–389, Jul 2018.



Y. Zhang, Y. Lau, H.-W. Kuo, S. Cheung, A. Pasupathy, and J. Wright.
On the global geometry of sphere-constrained sparse blind deconvolution.
In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.