



A Variable Smoothing for Nonconvexly Constrained Nonsmooth Optimization with Application to Sparse Spectral Clustering



東京工業大学
Tokyo Institute of Technology

Keita Kume and Isao Yamada, Tokyo Institute of Technology

Summary

Target problem:

$$\underset{x \in C}{\text{Minimize}} f(x) := \underbrace{h(x)}_{\text{smooth}} + \underbrace{g \circ G(x)}_{\text{nonsmooth}} \cdots (\star)$$

- \mathcal{X}, \mathcal{Z} are Euclidean spaces
- $C (\neq \emptyset) \subset \mathcal{X}$ is a (possibly nonconvex) closed subset of \mathcal{X}
- $h: \mathcal{X} \rightarrow \mathbb{R}$: differentiable, ∇h is Lipschitz continuous over C
- $g: \mathcal{Z} \rightarrow \mathbb{R}$: weakly convex, nonsmooth Lipschitz continuous
 $\Leftrightarrow \exists \eta > 0$ s. t. $g + \frac{\eta}{2} \|\cdot\|^2$ is convex
- $G: \mathcal{X} \rightarrow \mathcal{Z}$: smooth (possibly nonlinear) mapping

Typical applications: sparsity-aware signal processing, e.g., sparse PCA, sparse spectral clustering, robust subspace recovery

Challenging issues: $\left\{ \begin{array}{l} \text{nonconvex constraint } C \\ \text{nonsmoothness and nonconvexity of } g \end{array} \right\}$

Our contributions:

- Proposal of an optimization algorithm of guaranteed global convergence to a stationary point
(**First available algorithm for (\star) , and generalization of [1]**).
- Application to sparse spectral clustering (SSC) based on nonconvex sparse regularizer
(**Inherently first nonconvex approach for SSC**).

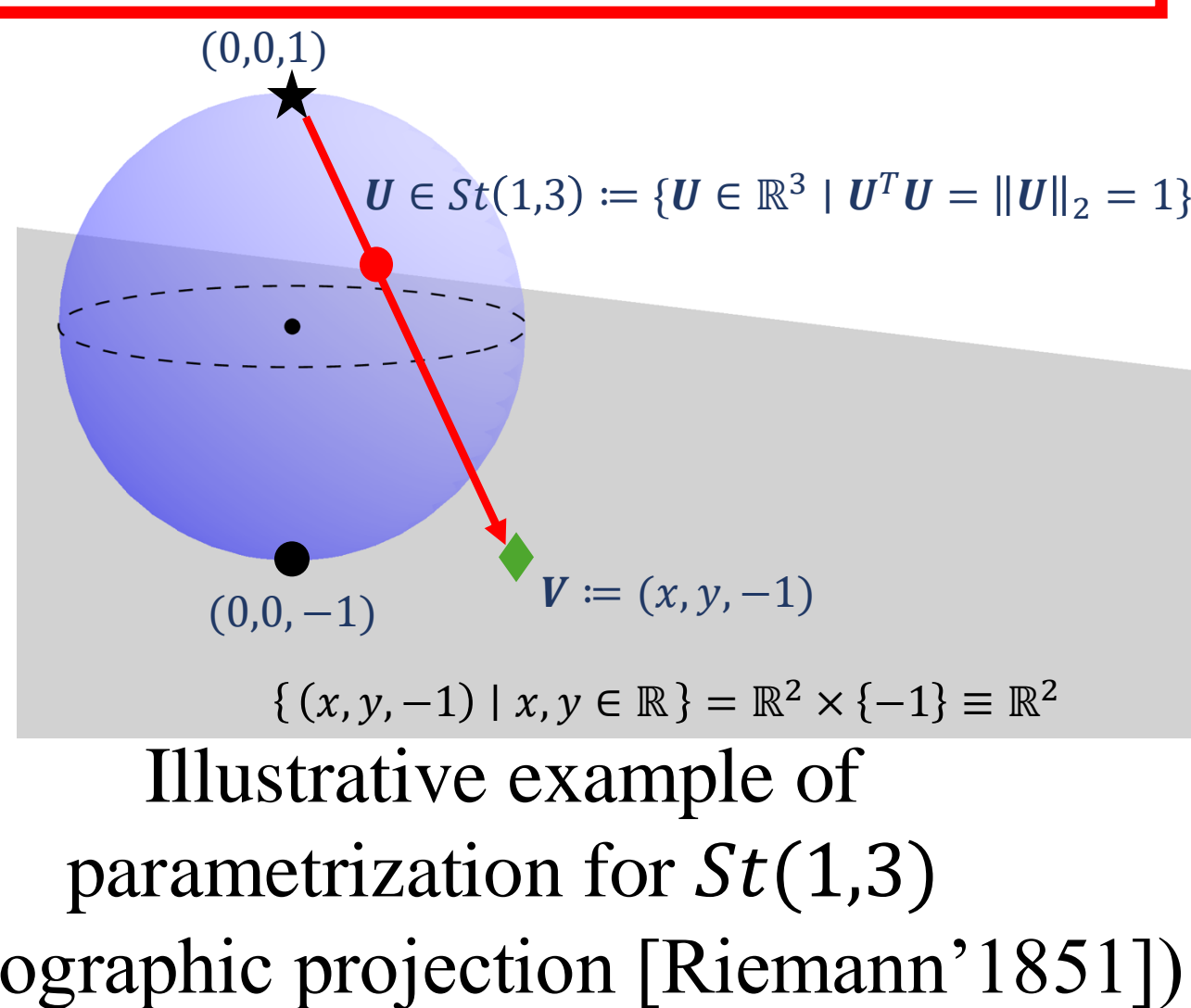
How to deal with constraint set C ? \Rightarrow parametrization

Key idea

Parameterize C in terms of the Euclidean space \mathcal{Y} with a smooth mapping $F: \mathcal{Y} \rightarrow \mathcal{X}$ such that $C = \{F(\mathbf{y}) \in \mathcal{X} \mid \mathbf{y} \in \mathcal{Y}\}$.

Example: smoothly parameterizable C

- Stiefel manifold
 $St(p, N) := \{U \in \mathbb{R}^{N \times p} \mid U^T U = I_p\}$
(with Cayley-type transforms [2,3])
- Bounded-rank matrices
 $\mathbb{R}_{\leq r}^{M \times N} := \{X \in \mathbb{R}^{M \times N} \mid \text{rank}(X) \leq r\}$
(with the multiplication $X = YZ^T$
 $[Y \in \mathbb{R}^{M \times r}, Z \in \mathbb{R}^{N \times r}]$ [4])



We consider the following parameterized problem instead of (\star) :

$$\underset{\mathbf{y} \in \mathcal{Y}}{\text{Minimize}} f \circ F(\mathbf{y}) = (h + g \circ G) \circ F(\mathbf{y}) \cdots (\clubsuit)$$

Euclidean space

We have the following relation of necessary conditions (optimality condition) of a local minimizer for (\star) and (\clubsuit) .

Theorem 4.1 [Relations of optimality conditions]

$$\mathbf{0} \in \partial f(F(\mathbf{y}^*)) + N_C(F(\mathbf{y}^*)) \Leftrightarrow \mathbf{0} \in \partial(f \circ F)(\mathbf{y}^*)$$

Optimality condition for (\star) Optimality condition for (\clubsuit)

under $\left\{ \begin{array}{l} \text{the Clarke regularity on } C \text{ (i.e., } C \text{ is sufficiently smooth)} \\ N_C(F(\mathbf{y}^*)) = \{x \in \mathcal{X} \mid (DF(\mathbf{y}^*))^*(x) = \mathbf{0}\} \text{ at } \mathbf{y}^* \in \mathcal{Y} \end{array} \right.$

$\star \partial f$ denotes the general subdifferential. Note: these are different senses from convex analysis (see [5]).
 $N_C(x) \subset \mathcal{X}$ denotes the general normal cone. $(DF(\mathbf{y}))^*$ denotes the adjoint of the Fréchet derivative (Jacobi matrix) at $\mathbf{y} \in \mathcal{Y}$.

How to deal with nonsmoothness of g ? \Rightarrow smoothing

Key idea (inspired by [1])

Use a smoothed surrogate function of η -weakly convex g .

\rightarrow Moreau envelope ${}^\mu g$ of g with $\mu \in (0, \eta^{-1})$

$$(\bar{z} \in \mathcal{Z}) \quad {}^\mu g(\bar{z}) := \inf_{z \in \mathcal{Z}} \left(g(z) + \frac{1}{2\mu} \|z - \bar{z}\|^2 \right)$$

- $\lim_{\mu \rightarrow 0} {}^\mu g(\bar{z}) = g(\bar{z})$
- ${}^\mu g$ is differentiable, and $\nabla {}^\mu g$ is Lipschitz continuous

Theorem 3.1 [Characterization of optimality condition]

For $(\mu_n)_{n=1}^\infty (\subset (0, \eta^{-1})) \xrightarrow{n \rightarrow \infty} 0$, and $(\mathbf{y}_n)_{n=1}^\infty \subset \mathcal{Y} \xrightarrow{n \rightarrow \infty} \exists \bar{\mathbf{y}} \in \mathcal{Y}$,
 $d(\mathbf{0}, \partial(f \circ F)(\bar{\mathbf{y}})) \leq \liminf_{n \rightarrow \infty} \|\nabla((h + \mu_n g \circ G) \circ F)(\mathbf{y}_n)\|$

$\liminf_{n \rightarrow \infty} \|\nabla((h + \mu_n g \circ G) \circ F)(\mathbf{y}_n)\| = 0$ implies
 $d(\mathbf{0}, \partial(f \circ F)(\bar{\mathbf{y}})) = 0$, i.e., $\mathbf{0} \in \partial(f \circ F)(\bar{\mathbf{y}})$.

Proposed algorithm achieves $\liminf_{n \rightarrow \infty} \|\nabla((h + \mu_n g \circ G) \circ F)(\mathbf{y}_n)\| = 0$:

- Set $\mu_n := \kappa n^{-\frac{1}{\alpha}}$ and $f_{[n]} := h + \mu_n g \circ G$ ($\alpha > 1, \exists \gamma_n > 0, \exists \kappa > 0$)
- Update $\mathbf{y}_{n+1} := \mathbf{y}_n - \gamma_n \nabla(f_{[n]} \circ F)(\mathbf{y}_n)$ Increment n

Theorem 3.3 [Convergence analysis (informal)]

Assume that $\nabla(f_{[n]} \circ F)$ is Lipschitz continuous with a Lipschitz constant $\varpi \mu_n^{-1}$ with some $\varpi > 0$, and $\gamma_n > 0$ is computed by the so-called *backtracking algorithm*. Then, $(\mathbf{y}_n)_{n=1}^\infty$ generated by the proposed algorithm satisfies:

$$\liminf_{n \rightarrow \infty} \|\nabla((h + \mu_n g \circ G) \circ F)(\mathbf{y}_n)\| = 0.$$

Application to sparse Spectral Clustering (SC)

Goal: split given data $(\xi_i)_{i=1}^N \subset \mathbb{R}^d$ into K groups without labeled data.

Outline of SC [6]

- Construct a similarity graph \mathcal{G} of $(\xi_i)_{i=1}^N$.
 $\left\{ \begin{array}{l} D \in \mathbb{R}^{N \times N}: \text{degree matrix} \\ W \in \mathbb{R}^{N \times N}: \text{adjacency matrix} \end{array} \right.$
- Compute K smallest eigenvectors $U^* \in St(K, N)$ of the graph Laplacian $L := I - D^{-\frac{1}{2}} W D^{-\frac{1}{2}} \in \mathbb{R}^{N \times N}$.
- Apply k-means algorithm to N row (normalized) vectors of U^* .

(Steps 2 and 3 correspond to splitting \mathcal{G} into K connected subgraphs)

To improve SC, the Sparse SC (SSC) utilizes a prior knowledge that $U^* U^{*T}$ is sparse (block diagonal) in the ideal case [7].

Step 2 of SC can be refined along SSC as:

$$\text{Find } U^* \in \underset{U \in St(K, N)}{\text{argmin}} \quad \underbrace{\text{Tr}(U^T L U)}_{\text{Eigenvalue decomp.}} + \underbrace{\lambda \psi(U U^T)}_{\text{Promote sparsity of } U^* U^{*T}} \cdots (\spadesuit)$$

($\psi: \mathbb{R}^{N \times N} \rightarrow \mathbb{R}$: sparsity promoting function)

(\spadesuit) can be reformulated as (\clubsuit) with $h(U) := \text{Tr}(U^T L U)$,
 $g := \lambda \psi$, $G(U) := U U^T$ and the generalized Cayley transform [3]

$$F := \Phi_S^{-1}: \mathcal{Y} \rightarrow St(p, N): Y \mapsto S(I - Y)(I + Y)^{-1} I_{N \times p},$$

$$\text{where } \mathcal{Y} := \left\{ \begin{bmatrix} A & -B^T \\ B & 0 \end{bmatrix} \in \mathbb{R}^{N \times N} \mid A^T = -A \in \mathbb{R}^{p \times p}, B \in \mathbb{R}^{(N-p) \times p} \right\}.$$

($St(p, N)$ and F above satisfy the assumption in Theorem 4.1)

We propose to solve (\spadesuit) with MCP (Minimax Concave Penalty) [8] as ψ .

Result: **the proposed SSC with MCP achieves the best performance!**

	iris		shuttle		segmentation		breast cancer		glass		wine		seeds	
	NMI	ARI	NMI	ARI	NMI	ARI	NMI	ARI	NMI	ARI	NMI	ARI	NMI	ARI
SC [6]	0.778	0.745	0.387	0.205	0.501	0.341	0.417	0.419	0.321	0.174	0.433	0.363	0.662	0.659
SSC(ℓ_1 +relax) [7]	0.785	0.786	0.426	0.279	0.503	0.343	0.433	0.462	0.325	0.176	0.433	0.363	0.671	0.675
Proposed SSC(ℓ_1 +Gr)	0.823	0.818	0.427	0.276	0.501	0.341	0.473	0.547	0.323	0.175	0.433	0.363	0.667	0.668
Proposed SSC(MCP+Gr)	0.823	0.818	0.434	0.294	0.507	0.351	0.558	0.664	0.341	0.180	0.442	0.376	0.721	0.708



arXiv version

[1] A. Böhm, S.J. Wright, "Variable smoothing for weakly convex composite functions," J. Optim. Theory Appl., 2021.
[2] K. Kume, I. Yamada, "A global Cayley parametrization of Stiefel manifold for direct utilization of optimization mechanisms over vector spaces," ICASSP, 2021.
[3] K. Kume, I. Yamada, "Generalized left-localized Cayley parametrization for optimization with orthogonality constraints," Optimization, 2022.
[4] E. Levin, J. Kileel, N. Boumal, "The effect of smooth parametrizations on nonconvex optimization landscapes," Math. Program., 2024.
[5] R. Rockafellar, R. J.-B. Wets, Variational Analysis, Springer Verlag, 3rd edition, 2010.
[6] A. Ng, M. Jordan, Y. Weiss, "On spectral clustering: Analysis and an algorithm," NeurIPS, 2001.
[7] C. Lu, S. Yan, Z. Lin, "Convex sparse spectral clustering: Single-view to multi-view," IEEE Trans. Image Process., 2016.
[8] C.-H. Zhang, "Nearly unbiased variable selection under mini-max concave penalty," The Annals of Statistics, 2010.