# Matrix Cones and Spectral Operators of Matrices

**Defeng Sun**

**Department of Applied Mathematics**

THE HONG KONG
POLYTECHNIC UNIVERSITY
香港理工大學

Advances in the Geometric and Analytic Theory of Convex Cones, May 29, 2019/Jeju

Based on joint works with Chao Ding, Jie Sun, and Kim-Chuan Toh

## The Metric Projector over the PSD Cone

1. Let $\mathcal{S}^n$ be set of $n$ by $n$ symmetric matrices in $\mathbb{R}^{m \times n}$ and $\mathcal{S}^n_+$ be the cone of positive semidefinite matrices in $\mathcal{S}^n$.

2. Let $X \in \mathcal{S}^n$ have the following spectral decomposition

$$X = P\Lambda P^{\mathbb{T}} = \sum_{i=1}^{n} \lambda_i p_i p_i^{\mathbb{T}},$$

where $\Lambda$ is the diagonal matrix of eigenvalues $\lambda_1, \ldots, \lambda_n$ of $X$ and $P$ is a corresponding orthogonal matrix of orthonormal eigenvectors. Then

$$X_+ := \Pi_{\mathcal{S}^n_+}(X) = P\Lambda_+ P^{\mathbb{T}} = \sum_{i=1}^{n} (\lambda_i)_+ p_i p_i^{\mathbb{T}}.$$

Here $\Pi_{\mathcal{S}^n_+}(X)$ is the unique optimal solution to

$$
\begin{aligned}
\min \quad & \frac{1}{2}\|Z - X\|_F^2 \\
\text{s.t.} \quad & Z \in \mathcal{S}^n_+ .
\end{aligned}
$$

1. Let $f : \Re \to \Re$ be a scalar function. The corresponding Löwner operator $F : \mathcal{S}^n \to \mathcal{S}^n$ is defined by[1]

$$F(X) := \sum_{i=1}^{n} f(\lambda_i) p_i p_i^{\mathbb{T}}, \quad X \in \mathcal{S}^n$$

2. Let $g : \Re \to \Re$ be an odd scalar function satisfying $g(-t) = -g(t)$ for all $t \geq 0$ (naturally $g(0) = 0$). One may define Löwner's operator $G : \mathbb{R}^{m \times n} \to \mathbb{R}^{m \times n}$ (assuming $m \leq n$) by

$$G(Z) := \sum_{i=1}^{m} g(\sigma_i(Z)) u_i v_i^{\mathbb{T}}, \quad Z \in \mathbb{R}^{m \times n},$$

where for any given $Z \in \mathbb{R}^{m \times n}$, $\sigma_1(Z) \geq \sigma_2(Z) \geq \ldots \geq \sigma_m(Z)$ denotes the singular values of $Z$ (always nonnegative and counting multiplicity) and $\sigma(Z)$ denotes the vector of the singular values of $Z$; $\mathbb{O}^{m,n}(Z)$ denotes the set of matrix pairs $(U, V) \in \mathbb{O}^m \times \mathbb{O}^n$ satisfying the singular value decomposition

$$Z = U \left[ \Sigma(Z) \quad 0 \right] V^{\mathbb{T}},$$

where $\Sigma(Z)$ is an $m \times m$ diagonal matrix whose $i$-th diagonal entry is $\sigma_i(Z) \geq 0$.

---

[1] Löwner, K.: *Über monotone matrixfunktionen*, Mathematische Zeitschrift 38 (1934) 177–216.

Let $X \in \mathbb{R}^{m \times n}$ admit the following singular value decomposition:

$$X = \overline{U} \left[ \Sigma(X) \ 0 \right] \overline{V}^T = \overline{U} \left[ \Sigma(X) \ 0 \right] \left[ \overline{V}_1 \ \overline{V}_2 \right]^T = \overline{U} \Sigma(X) \overline{V}_1^T, \tag{1}$$

where $\overline{U} \in \mathcal{O}^m$, $\overline{V} \in \mathcal{O}^n$ and $\overline{V}_1 \in \mathbb{R}^{n \times m}$, $\overline{V}_2 \in \mathbb{R}^{n \times (n-m)}$ and $\overline{V} = \left[ \overline{V}_1 \ \overline{V}_2 \right]$. The set of such matrices $(U, V)$ in the singular value decomposition (1) is denoted by $\mathcal{O}^{m,n}(X)$, i.e.,

$$\mathcal{O}^{m,n}(X) := \{ (U, V) \in \Re^{m \times m} \times \Re^{n \times n} \,|\, X = U \left[ \Sigma(X) \ 0 \right] V^T \}.$$

For any positive constant $\varepsilon > 0$, denote the closed convex cone $\mathcal{D}_n^\varepsilon$ by

$$\mathcal{D}_n^\varepsilon := \{(t,x) \in \mathbb{R} \times \mathbb{R}^n \,|\, \varepsilon^{-1}t \geq x_i, \ i = 1,\ldots,n\} \,. \tag{2}$$

Let $\Pi_{\mathcal{D}_n^\varepsilon}(\cdot)$ be the metric projector over $\mathcal{D}_n^\varepsilon$ under the Euclidean inner product in $\mathbb{R}^n$. That is, for any $(t,x) \in \mathbb{R} \times \mathbb{R}^n$, $\Pi_{\mathcal{D}_n^\varepsilon}(t,x)$ is the unique optimal solution to the following convex optimization problem

$$\begin{aligned} \min \quad & \frac{1}{2}\big((\tau - t)^2 + \|y - x\|^2\big) \\ \text{s.t.} \quad & \varepsilon^{-1}\tau \geq y_i, \ i = 1,\ldots,n \,. \end{aligned} \tag{3}$$

For any $x \in \mathbb{R}^n$, let $x^\downarrow$ be the vector of components of $x$ being arranged in the non-increasing order $x_1^\downarrow \geq \ldots \geq x_n^\downarrow$. Let $\mathrm{sgn}(x)$ be the sign vector of $x$, i.e., $(\mathrm{sgn})_i(x) = 1$ if $x_i \geq 0$ and $-1$ otherwise. We use " $\circ$ " to denote the Hadamard product operation either for two vectors or two matrices of the same dimensions.

### Proposition

*Assume that $\varepsilon > 0$ and $(t, x) \in \mathbb{R} \times \mathbb{R}^n$ are given. Let $\pi$ be a permutation of $\{1, \ldots, n\}$ such that $x^\downarrow = x_\pi$, i.e., $x_i^\downarrow = x_{\pi(i)}$, $i = 1, \ldots, n$ and $\pi^{-1}$ the inverse of $\pi$. For convenience, write $x_0^\downarrow = +\infty$ and $x_{n+1}^\downarrow = -\infty$. Let $\bar{\kappa}$ be the smallest integer $k \in \{0, 1, \ldots, n\}$ such that*

$$x_{k+1}^\downarrow \leq \Big( \sum_{j=1}^k x_j^\downarrow + \varepsilon t \Big) / (k + \varepsilon^2) < x_k^\downarrow. \tag{4}$$

Define $\bar{y} \in \mathbb{R}^n$ and $\bar{\tau} \in \mathbb{R}_+$, respectively, by

$$\bar{y}_i := \begin{cases} \Big( \sum_{j=1}^{\bar{\kappa}} x_j^{\downarrow} + \varepsilon t \Big) / (\bar{\kappa} + \varepsilon^2) & \text{if } 1 \le i \le \bar{\kappa}\,, \\ x_i^{\downarrow} & \text{otherwise} \end{cases}$$

and

$$\bar{\tau} := \varepsilon \bar{y}_1 = \varepsilon \Big( \sum_{j=1}^{\bar{k}} x_j^{\downarrow} + \varepsilon t \Big) / (\bar{k} + \varepsilon^2)\,.$$

The metric projection $\Pi_{\mathcal{D}_n^\varepsilon}(t, x)$ is computed by $\Pi_{\mathcal{D}_n^\varepsilon}(t, x) = (\bar{\tau}, \bar{y}_{\pi^{-1}})$.

For any positive constant $\varepsilon > 0$, define the matrix cone $\mathcal{M}_n^\varepsilon$ in $\mathcal{S}^n$ as the epigraph of the convex function $\varepsilon \lambda_{\max}(\cdot)$, i.e.,

$$\mathcal{M}_n^\varepsilon := \{(t, X) \in \mathbb{R} \times \mathcal{S}^n \mid \varepsilon^{-1} t \geq \lambda_{\max}(X)\}. \tag{5}$$

### Proposition

*Assume that $(t, X) \in \mathbb{R} \times \mathcal{S}^n$ is given. Let $X$ have the eigenvalue decomposition*

$$X = \overline{P} \operatorname{diag}(\lambda(X)) \overline{P}^T, \tag{6}$$

*where $\overline{P} \in \mathcal{O}^n$. Let $\Pi_{\mathcal{M}_n^\varepsilon}(\cdot, \cdot)$ be the metric projector over $\mathcal{M}_n^\varepsilon$ under Frobenius norm in $\mathcal{S}^n$. Then,*

$$\Pi_{\mathcal{M}_n^\varepsilon}(t, X) = (\bar{t}, \overline{P} \operatorname{diag}(\bar{y}) \overline{P}^T) \quad \forall\, (t, X) \in \mathbb{R} \times \mathcal{S}^n, \tag{7}$$

*where $(\bar{t}, \bar{y}) = \Pi_{\mathcal{D}_n^\varepsilon}(t, \lambda(X)) \in \Re \times \Re^n$.*

For any positive constant $\varepsilon > 0$, denote the closed convex cone $\mathcal{C}_n^\varepsilon$ by

$$\mathcal{C}_n^\varepsilon := \{(t, x) \in \mathbb{R} \times \mathbb{R}^n \,|\, \varepsilon^{-1} t \geq \|x\|_\infty\}. \tag{8}$$

Let $\Pi_{\mathcal{C}_n^\varepsilon}(\cdot, \cdot)$ be the metric projector over $\mathcal{C}_n^\varepsilon$ under the Euclidean inner product in $\mathbb{R}^n$. That is, for any $(t, x) \in \mathbb{R} \times \mathbb{R}^n$, $\Pi_{\mathcal{C}_n^\varepsilon}(t, x)$ is the unique optimal solution to the following convex optimization problem

$$\begin{aligned} \min \quad & \frac{1}{2}\big((\tau - t)^2 + \|y - x\|^2\big) \\ \text{s.t.} \quad & \varepsilon^{-1}\tau \geq \|y\|_\infty. \end{aligned} \tag{9}$$

In the following discussions, we frequently drop $n$ from $\mathcal{C}_n^\varepsilon$ when its size can be found from the context.

Assume that $\varepsilon > 0$ and $(t, x) \in \mathbb{R} \times \mathbb{R}^n$ are given. Let $\pi$ be a permutation of $\{1, \ldots, n\}$ such that $|x|^\downarrow = |x|_\pi$, i.e., $|x|_i^\downarrow = |x|_{\pi(i)}$, $i = 1, \ldots, n$ and $\pi^{-1}$ be the inverse of $\pi$. Let $|x|_0^\downarrow = +\infty$ and $|x|_{n+1}^\downarrow = 0$. Let $s_0 = 0$ and $s_k = \sum_{i=1}^k |x|_i^\downarrow$, $k = 1, \ldots, n+1$. Let $\overline{k}$ be the smallest integer $k \in \{0, 1, \ldots, n\}$ such that

$$|x|_{k+1}^\downarrow \leq (s_k + \varepsilon t)/(k + \varepsilon^2) < |x|_k^\downarrow \tag{10}$$

or $\overline{k} = n + 1$ if such an integer does not exist. Denote

$$\theta^\varepsilon(t, x) := (s_{\overline{k}} + \varepsilon t)/(\overline{k} + \varepsilon^2). \tag{11}$$

Let $\alpha, \beta$ and $\gamma$ be the index sets of $|x|^{\downarrow}$ as

$$\alpha := \{i \mid |x|_i^{\downarrow} > \theta^{\varepsilon}(t,x)\}, \quad \beta := \{i \mid |x|_i^{\downarrow} = \theta^{\varepsilon}(t,x)\} \tag{12}$$

and

$$\gamma := \{i \mid |x|_i^{\downarrow} < \theta^{\varepsilon}(t,x)\} \,. \tag{13}$$

Define $\bar{y} \in \mathbb{R}^n$ and $\bar{\tau} \in \Re_+$, respectively, by

$$\bar{y}_i := \left\{ \begin{array}{ll} \max\{\theta^{\varepsilon}(t,x), 0\} & \text{if } i \in \alpha \,, \\ |x|_i^{\downarrow} & \text{otherwise} \end{array} \right.$$

and

$$\bar{\tau} := \varepsilon \max\{\theta^{\varepsilon}(t,x), 0\} \,.$$

**Proposition**

Assume that $\varepsilon > 0$ and $(t, x) \in \mathbb{R} \times \mathbb{R}^n$ are given. The metric projection $\Pi_{\mathcal{C}^\varepsilon}(t, x)$ of $(t, x)$ onto $\mathcal{C}^\varepsilon$ can be computed as follows

$$\Pi_{\mathcal{C}^\varepsilon}(t, x) = \left( \bar{\tau}, \operatorname{sgn}(x) \circ \bar{y}_{\pi^{-1}} \right). \tag{14}$$

**Theorem**

Assume that $(t, X) \in \mathbb{R} \times \mathbb{R}^{m \times n}$ is given. Let $X$ have the singular value decomposition (1). Let $\Pi_{\mathcal{K}^\varepsilon}(\cdot, \cdot)$ be the metric projector over $\mathcal{K}^\varepsilon$ under Frobenius norm in $\mathbb{R}^{m \times n}$, where

$$\mathcal{K}^\varepsilon := \{(t, X) \in \mathbb{R} \times \mathbb{R}^{m \times n} \,|\, \varepsilon^{-1} t \geq \|X\|_2 \}. \tag{15}$$

For any $(t, X) \in \mathbb{R} \times \mathbb{R}^{m \times n}$, we have

$$\Pi_{\mathcal{K}^\varepsilon}(t, X) = \left( \bar{t}, \overline{U} \left[ \operatorname{diag}(\bar{y}) \ \ 0 \right] \overline{V}^T \right), \tag{16}$$

where

$$(\bar{t}, \bar{y}) = \Pi_{\mathcal{C}^\varepsilon}(t, \sigma(X)) \in \Re \times \Re^m.$$

12

## Matrix Optimization

1. Löwner operators are inadequate for applications

2. For a given unitarily invariant proper closed convex function $f : \mathcal{X} \to (-\infty, \infty]$, in matrix optimization one often considers the proximal mapping of $f$ at $X$:

$$\mathsf{P}_f(X) := \mathrm{argmin}_{Y \in \mathcal{X}} \left\{ f(Y) + \frac{1}{2} \|Y - X\|^2 \right\}, \quad X \in \mathcal{X}, \tag{17}$$

where $\mathcal{X}$ is either the real vector subspace $\mathbb{S}^m$ of $m \times m$ real symmetric (or complex) Hermitian matrices, or the real vector subspace $\mathbb{V}^{m \times n}$ of $m \times n$

3. For example, for $f(Y) = \|Y\|_2 = \sigma_{\max}(Y)$, the spectral norm of $Y$, $\mathsf{P}_f(\cdot)$ is no longer the Löwner operator [it is the Löwner operator for $f(Y) = \|Y\|_* = \sum_{i=1}^m \sigma_i(Y)$].

4. If $f(\cdot)$ is the indicator function of a matrix cone, then the proximal mapping $\mathsf{P}_f(\cdot)$ is the metric projector over the corresponding matrix cone.

Let $s$ be a positive integer and $0 \le s_0 \le s$ be a nonnegative integer. For given positive integers $m_1, \ldots, m_s$ and $n_{s_0+1}, \ldots, n_s$, define the real vector space $\mathcal{X}$ by

$$\mathcal{X} := \mathbb{S}^{m_1} \times \ldots \times \mathbb{S}^{m_{s_0}} \times \mathbb{V}^{m_{s_0+1} \times n_{s_0+1}} \times \ldots \times \mathbb{V}^{m_s \times n_s}. \tag{18}$$

Without loss of generality, we assume that $m_k \le n_k$, $k = s_0 + 1, \ldots, s$.

For any $X = (X_1, \ldots, X_s) \in \mathcal{X}$, we have for $1 \le k \le s_0$, $X_k \in \mathbb{S}^{m_k}$ and $s_0 + 1 \le k \le s$, $X_k \in \mathbb{V}^{m_k \times n_k}$. Denote

$$\mathcal{Y} := \mathbb{R}^{m_1} \times \ldots \times \mathbb{R}^{m_{s_0}} \times \mathbb{R}^{m_{s_0}} \times \ldots \times \mathbb{R}^{m_s}. \tag{19}$$

For any $X \in \mathcal{X}$, define $\kappa(X) \in \mathcal{Y}$ by

$$\kappa(X) := (\lambda(X_1), \ldots, \lambda(X_{s_0}), \sigma(X_{s_0+1}), \ldots, \sigma(X_s)).$$

Define the set $\mathcal{P}$ by

$$\mathcal{P} := \{(Q_1, \ldots, Q_s) \mid Q_k \in \mathbb{P}^{m_k}, \ 1 \le k \le s_0 \text{ and } Q_k \in \pm\mathbb{P}^{m_k}, \ s_0 + 1 \le k \le s\}.$$

Let $g : \mathcal{Y} \to \mathcal{Y}$ be a given mapping. For any $x = (x_1, \ldots, x_s) \in \mathcal{Y}$ with $x_k \in \mathbb{R}^{m_k}$, we write $g(x) \in \mathcal{Y}$ in the form $g(x) = (g_1(x), \ldots, g_s(x))$ with $g_k(x) \in \mathbb{R}^{m_k}$ for $1 \le k \le s$.

### Definition

The given mapping $g : \mathcal{Y} \to \mathcal{Y}$ is said to be *mixed symmetric*, with respect to $\mathcal{P}$, at $x = (x_1, \ldots, x_s) \in \mathcal{Y}$ with $x_k \in \mathbb{R}^{m_k}$, if

$$g(Q_1 x_1, \ldots, Q_s x_s) = (Q_1 g_1(x), \ldots, Q_s g_s(x)) \quad \forall \, (Q_1, \ldots, Q_s) \in \mathcal{P}. \tag{20}$$

The mapping $g$ is said to be mixed symmetric, with respect to $\mathcal{P}$, over a set $\mathcal{D} \subseteq \mathcal{Y}$ if (20) holds for every $x \in \mathcal{D}$. We call $g$ a *mixed symmetric* mapping, with respect to $\mathcal{P}$, if (20) holds for every $x \in \mathcal{Y}$.

## Spectral Operators

Note that for each $k \in \{1, \ldots, s\}$, the function value $g_k(x) \in \mathbb{R}^{m_k}$ is dependent on all $x_1, \ldots, x_s$. When there is no danger of confusion, in later discussions we often drop the phrase "with respect to $\mathcal{P}$" from Definition 1. Let $\mathcal{N}$ be a given nonempty set in $\mathcal{X}$. Define $\kappa_{\mathcal{N}} := \{\kappa(X) \in \mathcal{Y} \mid X \in \mathcal{N}\}$. The following definition of the spectral operator with respect to a mixed symmetric mapping $g$.

### Definition

Suppose that $g : \mathcal{Y} \to \mathcal{Y}$ is mixed symmetric on $\kappa_{\mathcal{N}}$. The spectral operator $G : \mathcal{N} \to \mathcal{X}$ with respect to $g$ is defined as $G(X) := (G_1(X), \ldots, G_s(X))$ for $X = (X_1, \ldots, X_s) \in \mathcal{N}$ such that

$$G_k(X) := \begin{cases} P_k \mathrm{Diag}\big(g_k(\kappa(X))\big) P_k^{\mathbb{T}} & \text{if } 1 \leq k \leq s_0, \\ U_k \begin{bmatrix} \mathrm{Diag}\big(g_k(\kappa(X))\big) & 0 \end{bmatrix} V_k^{\mathbb{T}} & \text{if } s_0 + 1 \leq k \leq s, \end{cases}$$

where $P_k \in \mathbb{O}^{m_k}(X_k)$, $1 \leq k \leq s_0$, $(U_k, V_k) \in \mathbb{O}^{m_k, n_k}(X_k)$, $s_0 + 1 \leq k \leq s$.

## Spectral Operators

Next, we will focus on the study of spectral operators for the case that $\mathcal{X} \equiv \mathbb{V}^{m \times n}$. The corresponding extensions for the spectral operators defined on the general Cartesian product of several matrix spaces can be considered in a similar fashion.

Let $\mathcal{N}$ be a given nonempty open set in $\mathbb{V}^{m \times n}$. Suppose that $g : \mathbb{R}^m \to \mathbb{R}^m$ is mixed symmetric with respect to $\mathcal{P} \equiv \pm \mathbb{P}^m$ (i.e., absolutely symmetric), on an open set $\widehat{\sigma}_{\mathcal{N}}$ in $\mathbb{R}^m$ containing $\sigma_{\mathcal{N}} := \{\sigma(X) \mid X \in \mathcal{N}\}$. The spectral operator $G : \mathcal{N} \to \mathbb{V}^{m \times n}$ with respect to $g$ then takes the form of

$$G(X) = U \left[\mathrm{Diag}(g(\sigma(X))) \quad 0\right] V^{\mathbb{T}}, \quad X \in \mathcal{N},$$

where $(U, V) \in \mathbb{O}^{m,n}(X)$. For a given $\overline{X} \in \mathcal{N}$, consider the singular value decomposition (SVD) of $\overline{X}$, i.e.,

$$\overline{X} = \overline{U} \left[\Sigma(\overline{X}) \quad 0\right] \overline{V}^{\mathbb{T}}, \tag{21}$$

where $\Sigma(\overline{X})$ is an $m \times m$ diagonal matrix whose $i$-th diagonal entry is $\sigma_i(\overline{X})$, $\overline{U} \in \mathbb{O}^m$ and $\overline{V} = \left[\overline{V}_1 \quad \overline{V}_2\right] \in \mathbb{O}^n$ with $\overline{V}_1 \in \mathbb{V}^{n \times m}$ and $\overline{V}_2 \in \mathbb{V}^{n \times (n-m)}$.

1. Let $\mathcal{X}, \mathcal{Y}$ be two finite-dimensional real Euclidean spaces
2. $F : \mathcal{X} \to \mathcal{Y}$ a locally Lipschitz continuous function.

Since $F$ is almost everywhere differentiable [Rademacher, 1912], we can define

$$\partial_B F(x) := \left\{\lim F'(x^k) : \ x^k \to x, \ x^k \in D_F\right\}.$$

Here $D_F$ is the set of points where $F$ is differentiable. Hence, Clarke's generalized Jacobian of $F$ at $x$ is given by

$$\partial F(x) = \operatorname{conv} \partial_B F(x).$$

### Definition

Let $\mathcal{K} : \mathcal{X} \rightrightarrows \mathcal{L}(\mathcal{X}, \mathcal{Y})$ be a nonempty, compact valued and upper-semicontinous multifunction. We say that $F$ is semismooth $x \in \mathcal{X}$ with respect to $\mathcal{K}$ if (i) $F$ is directionally differentiable at $x$; and (ii) for any $\Delta x \in \mathcal{X}$ and $V \in \mathcal{K}(x + \Delta x)$ with $\Delta x \to 0$,

$$F(x + \Delta x) - F(x) - V(\Delta x) = o(\|\Delta x\|) \quad (g-semismooth). \tag{22}$$

Furthermore, if (22) is replaced by

$$F(x + \Delta x) - F(x) - V(\Delta x) = O(\|\Delta x\|^{1+\gamma}), \tag{23}$$

where $\gamma > 0$ is a constant, then $F$ is said to be $\gamma$-order (strongly if $\gamma = 1$) semismooth at $x$ with respect to $\mathcal{K}$.

Assume that $F(\bar{x}) = 0$.

Given $x^0 \in \mathcal{X}$. For $k = 0, 1, \ldots$

Main Step   Choose an arbitrary $V_k \in \mathcal{K}(x^k)$. Solve

$$F(x^k) + V_k(x^{k+1} - x^k) = 0$$

Rates of Convergence: Assume that $\mathcal{K}(\bar{x})$ is nonsingular and that $x^0$ is sufficiently close to $\bar{x}$. If $F$ is g-semismooth at $\bar{x}$, then

$$\|x^{k+1} - \bar{x}\| = \| \underbrace{V_k^{-1}}_{\text{bounded}} \underbrace{[F(x^k) - F(\bar{x}) - V_k(x^k - \bar{x})]}_{\text{g-semismooth}}\| = \underbrace{o(\|x^k - \bar{x}\|)}_{\text{superlinear}}.$$

It takes $o(\|x^k - \bar{x}\|^{1+\gamma})$ if $F$ is $\gamma$-order g-semismooth at $\bar{x}$ [the directional differentiability of $F$ is not needed in the above local convergence analysis]

## Nonsmooth Equations

1. The nonsmooth equation approach is popular in the complementarity and variational inequalities (nonsmooth equations) community (Robinson, Pang, ...)

2. Josephy (1979) introduced Newton and quasi-Newton methods for generalized equations (in terms of Robinson).

3. Kojima and Shindo (1986) investigated Newton's method for piecewise smooth equations.

4. Kummer (1988, 1992) gave a sufficient condition (22) to extend Kojima and Shindo's work.

5. L. Qi and J. Sun (1993) proved what we know now.

6. Since then, many exciting developments, in particular in the large-scale settings ...

Why nonsmooth Newton methods important in solving large-scale optimization problems? We illustrate this with an example.

## The nearest correlation matrix problem: An example

Consider the nearest correlation matrix (NCM) problem:

$$\min \left\{ \frac{1}{2} \| X - G \|_F^2 \mid X \succeq 0, \ X_{ii} = 1, \ i = 1, \ldots, n \right\}.$$

The dual of the above problem can be written as (in its minimization format)

$$\min \quad \frac{1}{2} \| \Xi \|^2 - \langle b, \ y \rangle - \frac{1}{2} \| G \|^2$$
$$\text{s.t.} \quad S - \Xi + \mathcal{A}^* y = -G, \quad S \succeq 0$$

or via eliminating $\Xi$ and $S \succeq 0$, the following

$$\min \left\{ \varphi(y) := \frac{1}{2} \| \Pi_{\mathcal{S}_+^n} (\mathcal{A}^* y + G) \|^2 - \langle b, \ y \rangle - \frac{1}{2} \| G \|^2 \right\},$$

which is equivalent to the strongly semismooth system (S. & Sun, 02) of equations

$$\nabla \varphi(y) = \mathcal{A} \Pi_{\mathcal{S}_+^n} (\mathcal{A}^* y + G) - b = 0.$$

Test the second order nonsmooth Newton-CG method [H.-D. Qi & S. 06] ([X,y] = CorrelationMatrix(G,b,tau,tol) in Matlab from Sun's webpage) and two popular first order methods (FOMs) [APG of Nesterov; ADMM of Glowinski (steplength $1.618$)] all to the dual forms for the NCM with real financial data:

$G$: Cor3120, $n = 3,120$, obtained from [N. J. Higham & N. Strabić, SIMAX, 2016] [Optimal sol. rank $= 3,025$, high rank]

| $n = 3,120$ | Newton-CG | ADMM | APG |
|---|---|---|---|
| Rel. KKT Res. | 2.7-8 | 2.9-7 | 9.2-7 |
| time (s) | 26.8 | 246.4 | 459.1 |
| iters | 4 | 58 | 111 |
| avg-time/iter | 6.7 | 4.3 | 4.1 |

Newton's method only takes at most $40\%$ time more than ADMM & APG (or FISTA) per iteration (Newton will take less time on average per iteration if it took more iterations).

### Theorem

*Suppose that $\overline{X} \in \mathcal{N}$ has the SVD (21). The spectral operator $G$ is continuous at $\overline{X}$ if and only if $g$ is continuous at $\sigma(\overline{X})$.*

### Theorem

*Suppose that $\overline{X}$ has the SVD (21). The spectral operator $G$ is locally Lipschitz continuous near $\overline{X}$ if and only if $g$ is locally Lipschitz continuous near $\overline{\sigma} = \sigma(\overline{X})$.*

Let $\eta(\sigma) \in \mathbb{R}^m$ be the vector defined by ($i \in \{1, \ldots, m\}$)

$$(\eta(\sigma))_i := \left\{ \begin{array}{ll} (g'(\sigma))_{ii} - (g'(\sigma))_{ij} & \text{if } \exists j \in \{1, \ldots, m\} \text{ and } j \neq i \text{ such that } \sigma_i = \sigma_j, \\ (g'(\sigma))_{ii} & \text{otherwise}, \end{array} \right. \quad (24)$$

Define the corresponding *divided difference matrix* $\mathcal{E}_1(\sigma) \in \mathbb{R}^{m \times m}$, the *divided addition matrix* $\mathcal{E}_2(\sigma) \in \mathbb{R}^{m \times m}$, the *division matrix* $\mathcal{F}(\sigma) \in \mathbb{R}^{m \times (n-m)}$, respectively, by

$$(\mathcal{E}_1(\sigma))_{ij} := \left\{ \begin{array}{ll} \dfrac{g_i(\sigma) - g_j(\sigma)}{\sigma_i - \sigma_j} & \text{if } \sigma_i \neq \sigma_j, \\ (\eta(\sigma))_i & \text{otherwise}, \end{array} \right. \quad i, j \in \{1, \ldots, m\}, \quad (25)$$

$$(\mathcal{E}_2(\sigma))_{ij} := \left\{ \begin{array}{ll} \dfrac{g_i(\sigma) + g_j(\sigma)}{\sigma_i + \sigma_j} & \text{if } \sigma_i + \sigma_j \neq 0, \\ (g'(\sigma))_{ii} & \text{otherwise}, \end{array} \right. \quad i, j \in \{1, \ldots, m\}, \quad (26)$$

$$(\mathcal{F}(\sigma))_{ij} := \left\{ \begin{array}{ll} \dfrac{g_i(\sigma)}{\sigma_i} & \text{if } \sigma_i \neq 0, \\ (g'(\sigma))_{ii} & \text{otherwise}, \end{array} \right. \quad i \in \{1, \ldots, m\}, \quad j \in \{1, \ldots, n-m\}. \quad (27)$$

Define the matrix $\mathcal{C}(\sigma) \in \mathbb{R}^{m \times m}$ to be the difference between $g'(\sigma)$ and $\mathrm{Diag}(\eta(\sigma))$, i.e.,

$$\mathcal{C}(\sigma) := g'(\sigma) - \mathrm{Diag}(\eta(\sigma)). \tag{28}$$

When the dependence of $\eta$, $\mathcal{E}_1$, $\mathcal{E}_2$, $\mathcal{F}$ and $\mathcal{C}$ on $\sigma$ is clear from the context, we often drop $\sigma$ from the corresponding notations. Note that the divided difference matrix $\mathcal{E}_1(\sigma)$ is the same with the commonly defined for the symmetric matrix case. The divided addition matrix $\mathcal{E}_2(\sigma)$ and the division matrix $\mathcal{F}(\sigma)$ are particular to general non-Hermitian matrices.

Denote $\overline{\eta} = \eta(\overline{\sigma}) \in \mathbb{R}^m$ to be the vector defined by (24). Let $\overline{\mathcal{E}}_1$, $\overline{\mathcal{E}}_2$, $\overline{\mathcal{F}}$ and $\overline{\mathcal{C}}$ be the real matrices defined in (25)–(28) with respect to $\overline{\sigma}$.

### Theorem

*Suppose that the given matrix $\overline{X} \in \mathcal{N}$ has the SVD (21). Then the spectral operator $G$ is F-differentiable at $\overline{X}$ if and only if $g$ is F-differentiable at $\overline{\sigma}$. In that case, the derivative of $G$ at $\overline{X}$ is given by*

$$G'(\overline{X})H = \overline{U}[\overline{\mathcal{E}}_1 \circ S(A) + \mathrm{Diag}\left(\overline{\mathcal{C}}\mathrm{diag}(S(A))\right) + \overline{\mathcal{E}}_2 \circ T(A) \quad \overline{\mathcal{F}} \circ B]\overline{V}^{\mathbb{T}} \quad \forall H \in \mathbb{V}^{m \times n}, \text{ (29)}$$

*where $A := \overline{U}^{\mathbb{T}}H\overline{V}_1$, $B := \overline{U}^{\mathbb{T}}H\overline{V}_2$ and for any $X \in \mathbb{V}^{m \times m}$, $\mathrm{diag}(X)$ denotes the column vector consisting of all the diagonal entries of $X$ being arranged from the first to the last.*

Here the two linear matrix operators $S : \mathbb{V}^{p \times p} \to \mathbb{S}^p$ and $T : \mathbb{V}^{p \times p} \to \mathbb{V}^{p \times p}$ are given by

$$S(Y) := \frac{1}{2}(Y + Y^{\mathbb{T}}), \quad T(Y) := \frac{1}{2}(Y - Y^{\mathbb{T}}), \quad Y \in \mathbb{V}^{p \times p}. \tag{30}$$

## B(ouligand)-Differentiability

Let $\mathcal{Z}$ be a finite dimensional real Euclidean space equipped with an inner product $\langle \cdot, \cdot \rangle$ and its induced norm $\| \cdot \|$. Let $\mathcal{O}$ be an open set in $\mathcal{Z}$ and $\mathcal{Z}'$ be another finite dimensional real Euclidean space. The function $F : \mathcal{O} \subseteq \mathcal{Z} \to \mathcal{Z}'$ is said to be *B(ouligand)-differentiable* at $z \in \mathcal{O}$ if for any $h \in \mathcal{Z}$ with $h \to 0$,

$$F(z + h) - F(z) - F'(z; h) = o(\|h\|).$$

A stronger notion than B-differentiability is $\rho$-order B-differentiability with $\rho > 0$. The function $F : \mathcal{O} \subseteq \mathcal{Z} \to \mathcal{Z}'$ is said to be *$\rho$-order B-differentiable* at $z \in \mathcal{O}$ if for any $h \in \mathcal{Z}$ with $h \to 0$,

$$F(z + h) - F(z) - F'(z; h) = O(\|h\|^{1+\rho}).$$

### Theorem

*Suppose that $\overline{X} \in \mathcal{N}$ has the SVD (21). Let $0 < \rho \leq 1$ be given.*

(i) *If $g$ is locally Lipschitz continuous near $\sigma(\overline{X})$ and $\rho$-order B-differentiable at $\sigma(\overline{X})$, then $G$ is $\rho$-order B-differentiable at $\overline{X}$.*

(ii) *If $G$ is $\rho$-order B-differentiable at $\overline{X}$, then $g$ is $\rho$-order B-differentiable at $\sigma(\overline{X})$.*

### Theorem

*Suppose that $\overline{X} \in \mathcal{N}$ has the singular value decomposition (21). Let $0 < \rho \leq 1$ be given. $G$ is $\rho$-order g-semismooth at $\overline{X}$ if and only if $g$ is $\rho$-order g-semismooth at $\overline{\sigma}$.*

## Characterizations of the Generalized Jacobians

Assume that $g$ is locally Lispchitz continuous. Then since the spectral operator $G$ is locally Lipschitz continuous near $\overline{X}$, $\Psi = G'(\overline{X}; \cdot)$ is globally Lipschitz continuous if exists. In that case, $\partial_B \Psi(0)$ and $\partial \Psi(0)$ are well-defined. Furthermore, we have the following characterization of the B-subdifferential and Clarke's subdifferential of the spectral operator $G$ at $\overline{X}$.

### Theorem

*Suppose that the given $\overline{X} \in \mathcal{N}$ has the decomposition (21). Suppose that there exists an open neighborhood $\mathcal{B} \subseteq \mathbb{R}^m$ of $\overline{\sigma}$ in $\widehat{\sigma}_{\mathcal{N}}$ such that $g$ is differentiable at $\sigma \in \mathcal{B}$ if and only if $g'(\overline{\sigma}; \cdot)$ is differentiable at $\sigma - \overline{\sigma}$. Assume further that the function $d : \mathbb{R}^m \to \mathbb{R}^m$ defined by*

$$d(h) := g(\overline{\sigma} + h) - g(\overline{\sigma}) - g'(\overline{\sigma}; h), \quad h \in \mathbb{R}^m \tag{31}$$

*is strictly differentiable at zero. Then, we have*

$$\partial_B G(\overline{X}) = \partial_B \Psi(0) \quad \text{and} \quad \partial G(\overline{X}) = \partial \Psi(0).$$

Many more to be developed ...

## References

1. Chao Ding, Defeng Sun, and Kim-Chuan Toh, "An introduction to a class of matrix cone programming," Mathematical Programming 144 (2014) 141-179.

2. Chao Ding, Defeng Sun, Jie Sun, and Kim-Chuan Toh, "Spectral operators of matrices," Mathematical Programming 168 (2018) 509–531.

3. Chao Ding, Defeng Sun, Jie Sun, and Kim-Chuan Toh, "Spectral operators of matrices: semismoothness and characterizations of the generalized Jacobian," arXiv:1810.09856, October 2018.