

# Several observations about using the ALM + semismooth Newton method for solving large scale SDP and beyond

**Defeng Sun**

**Department of Applied Mathematics**



THE HONG KONG  
POLYTECHNIC UNIVERSITY  
香港理工大學

One World Optimization Seminar  
November 30, 2020

# The nearest correlation matrix model

Let us start with a simple Nearest Correlation Matrix (NCM) problem (a terminology coined by Nick Higham in 2002)

$$\begin{aligned} \min \quad & \frac{1}{2} \|X - G\|_F^2 \\ \text{s.t.} \quad & X_{ii} = 1, \quad i = 1, \dots, n, \\ & X \succeq 0, \end{aligned}$$

which is a special case of the **best approximation problem**

$$\begin{aligned} \min \quad & \frac{1}{2} \|x - c\|^2 \\ \text{s.t.} \quad & \mathcal{A}x \in b + Q, \\ & x \in K, \end{aligned}$$

where  $\mathcal{X}$  is a finite-dimensional real Hilbert space equipped with a scalar product  $\langle \cdot, \cdot \rangle$  and its induced norm  $\|\cdot\|$

$\mathcal{A} : \mathcal{X} \rightarrow \mathbb{R}^m$  is a linear operator

$Q = \{0\}^p \times \mathbb{R}_+^q$  is a polyhedral convex cone,  $1 \leq p \leq m$ ,  $q = m - p$ , and  $K$  is a closed convex cone in  $\mathcal{X}$ .

The dual of the best approximation problem is in the form of

$$\begin{aligned} \max \quad & -\theta(y) := - \left[ \frac{1}{2} \|\Pi_K(c + \mathcal{A}^*y)\|^2 - \langle b, y \rangle - \frac{1}{2} \|c\|^2 \right] \\ \text{s.t.} \quad & y \in Q^*, \end{aligned}$$

where  $Q^* = \Re^p \times \Re_+^q$ .

One can solve the above dual by solving the following nonsmooth equation

$$F(y) := y - \Pi_{Q^*}[y - (\mathcal{A}\Pi_K(c + \mathcal{A}^*y) - b)] = 0, \quad y \in \Re^m.$$

For the **NCM problem**,

- $\mathcal{A}(X) = \text{diag}(X)$ , a vector consisting of all diagonal entries of  $X$ .
- $\mathcal{A}^*(y) = \text{diag}(y)$ , the diagonal matrix.
- $b = e$ , the vector of all ones in  $\Re^n$  and  $K = \mathcal{S}_+^n$

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

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

$$X = P\Lambda P^T,$$

where  $\Lambda$  is the diagonal matrix of eigenvalues of  $X$  arranged in the non-increasing order and  $P$  is a corresponding orthogonal matrix of orthonormal eigenvectors.

Then

$$X_+ := \Pi_{\mathcal{S}_+^n}(X) = P\Lambda_+P^T.$$

When  $X$  is nonsingular,  $\Pi_{\mathcal{S}_+^n}(\cdot)$  is continuously differentiable near  $X$  with

$$\Pi'_{\mathcal{S}_+^n}(X)(\Delta X) = P[\Omega \circ P^T(\Delta X)P]P^T \quad \forall \Delta X \in \mathcal{S}^n,$$

where

$$\Omega_{ij} = \frac{\max(0, \lambda_i) - \max(0, \lambda_j)}{\lambda_i - \lambda_j} \text{ if } \lambda_i \neq \lambda_j \text{ \& } \Omega_{ij} = (\max(0, \lambda_i))' \text{ otherwise.}$$

We have

- $\|X_+\|^2$  is continuously differentiable with

$$\nabla\left(\frac{1}{2}\|X_+\|^2\right) = X_+,$$

but is not twice continuously differentiable (good news!).

- $X_+$  is not piecewise smooth unless  $n = 1$ , but strongly semismooth<sup>1</sup>: for any  $X \in \mathcal{S}^n$ ,  $\Pi_{\mathcal{S}_+^n}$  is directionally differentiable at  $X$  and it holds for  $X$  that

$$\Pi_{\mathcal{S}_+^n}(Y) - \Pi_{\mathcal{S}_+^n}(X) - \Pi'_{\mathcal{S}_+^n}(Y)(Y - X) = O(\|Y - X\|^2)$$

for any  $Y \in \mathcal{S}^n$  such that  $\Pi_{\mathcal{S}_+^n}(\cdot)$  is differentiable at  $Y$  and  $Y \rightarrow X$ .

---

<sup>1</sup>S. and J. Sun, Semismooth matrix valued functions. *Math of OR* 27 (2002) 150–169.

We test the efficiency of a quadratically convergent Newton's method designed by Houduo Qi and S.<sup>2</sup> The written code is called CorrelationMatrix.m on randomly generated pseudo correlation matrix. The code (in other languages too) is publically available in my webpage. All the experiments are done on the ThinkStation Desktop with Intel (R) Core(TM) i7-8700 processor.

n	relgap		Iterations		Time (s)	
1000	1.70E-08	8.05E-14	5	7	0.9	1.2
3000	3.83E-08	6.41E-14	5	7	13.7	18.5
5000	8.68E-07	1.91E-12	5	7	61.79	83.8
10000	1.81E-07	1.95E-12	5	7	342.42	520.25

**Observations:** 1. In addition to the strong convexity of the objective function, the key point for the semismooth Newton method to work is the primal constraint non-degeneracy (transversality, generalized LICQ). 2.  $\Omega$  can be very sparse to reduce computational costs substantially.

---

<sup>2</sup>H.D. Qi and S., A quadratically convergent Newton method for computing the nearest correlation matrix. *SIAM J. Matrix Analysis and Applications* 28 (2006) 360–385.

If we have lower and upper bounds on  $X$ ,  $F$  takes the form

$$F(y) = y - \Pi_{Q^*} [y - (\mathcal{A}\Pi_{\mathcal{S}_+^n}(G + \mathcal{A}^*y) - b)],$$

which involves double layered projections over convex cones.

A quadratically convergent smoothing Newton method is designed by Yan Gao and S.<sup>3</sup>

Again, highly efficient if the constraint non-degeneracy holds!

---

<sup>3</sup>Y. Gao and S., Calibrating least squares covariance matrix problems with equality and inequality constraints, SIAM J on Matrix Analysis and Applications 31 (2009), 1432–1457.

One may write the NCM as a symmetric cone programming with both SDP cone and SOC cone constraints:

$$\begin{aligned} \min \quad & t \\ \text{s.t.} \quad & X_{ii} = 1, \quad i = 1, \dots, n, \\ & y + \text{svec}(X) = \text{svec}(G), \\ & X \in \mathcal{S}_+^n, \quad t \geq \|y\|_2. \end{aligned}$$

Then, we can solve the problem with interior point methods (IPMs).



We compare our algorithm with the state-of-the-art IPM based software SDPT3<sup>4</sup>.

n	relgap		Iterations		Time(s)	
	SDPT3	Newton	SDPT3	Newton	SDPT3	Newton
100	1.79E-09	7.80E-10	13	5	20.95	<b>0.02</b>
150	1.74E-06	1.74E-09	15	5	126.3	<b>0.03</b>
200	4.41E-09	2.54E-09	14	5	425.3	<b>0.04</b>
500	<b>X</b>	8.81E-09	<b>X</b>	5	<b>X</b>	<b>0.19</b>

**x** means SDPT3 is out of memory (117.8 GB memory is required).

---

<sup>4</sup>R.H TUTUNCU, K.C. TOH, AND M.J. TODD, Solving semidefinite-quadratic-linear programs using SDPT3, Mathematical Programming Ser. B, 95 (2003), pp. 189-217.

The dual problem of the NCM is an unconstrained optimization problem:

$$\max \quad -\theta(y) := -\left[\frac{1}{2}\|\Pi_{\mathcal{S}_+^n}(G + \mathcal{A}^*y)\|_F^2 - \langle b, y \rangle - \frac{1}{2}\|G\|_F^2\right],$$

with

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

Then, it is natural to apply the accelerated proximal gradient (APG) method to solve this dual problem:

$$\begin{cases} z^{k+1} = x^k - \nabla\theta(y^k), \\ x^{k+1} = \left(1 - \frac{2}{k+2}\right)x^k + \frac{2}{k+2}z^k, \\ y^{k+1} = \left(1 - \frac{2}{(k+1)+2}\right)x^{k+1} + \frac{2}{(k+1)+2}z^{k+1}, \end{cases}$$

where  $x^0 = z^0 = y^0$  and  $y^0$  is the given initial point.

We compare the APG with the semismooth Newton method on a randomly generated data, for the given integer  $n$ , set

$$G := \text{rand}(n, n), G = 0.5 * (G + G') - \text{diag}(\text{diag}(0.5 * (G + G'))) + \text{eye}(n).$$

We set the maximum iterations for APG to be 1,000 and the tolerance to be  $1e - 6$ .

n	Residual		Iterations		Time(s)	
	APG	Newton	APG	Newton	APG	Newton
100	9.62E-07	1.82E-07	53	4	0.08	<b>0.02</b>
500	9.95E-07	1.05E-07	104	5	1.5	<b>0.26</b>
1000	9.79E-07	4.22E-07	142	5	8.64	<b>0.8</b>
2000	9.48E-07	4.60E-07	191	5	93.01	<b>4.4</b>

To see the robustness of APG, we also test

$$G := \text{rand}(1000, 1000), G = G + G' - \text{diag}(\text{diag}(G + G')) + \text{eye}(1000).$$

We set the maximum iterations for APG to be 1,000 and the tolerance to be  $1e - 6$ .

Algorithms	Time (s)	Iterations	Residual
APG	61.1	1000	3.98E-04
Newton	1.4	9	3.34E-07

This motivates us to use semismooth Newton methods to solve semidefinite programming (SDP) under the framework of proximal point algorithms (PPAs).

Now we can turn to the following linear conic programming

$$\begin{aligned} \min \quad & \langle c, x \rangle \\ \text{s.t.} \quad & \mathcal{A}x = b, \\ & x \in \mathcal{K} \end{aligned}$$

where  $\mathcal{A} : \mathcal{X} \rightarrow \mathbb{R}^m$ ,  $b \in \mathbb{R}^m$ ,  $c \in \mathcal{X}$ ,  $\mathcal{X}$  is a finite-dimensional real Hilbert space and  $\mathcal{K} \subseteq \mathcal{X}$  is a closed convex cone, e.g., the second-order-cone or the PSD (positive and semidefinite) cone. Define the Lagrange function as follows:

$$L(x; y, s) := \langle c, x \rangle + \langle y, b - \mathcal{A}x \rangle - \langle s, x \rangle.$$

Let  $\mathcal{K}^*$  be the dual cone of  $\mathcal{K}$ . Then the (Lagrange) dual of linear conic programming is defined as

$$\max_{y \in \mathbb{R}^m, s \in \mathcal{K}^*} \left\{ \inf_{x \in \mathcal{X}} L(x; y, s) \right\}.$$

# Lagrange dual of linear conic programming

By noting that for any  $(x, y, s) \in \mathcal{X} \times \mathbb{R}^m \times \mathcal{X}^n$ ,

$$L(x; y, s) := \langle x, c - \mathcal{A}^*y - s \rangle + \langle y, b \rangle$$

we get an explicit formula for the dual problem as in the following

$$\begin{aligned} \max \quad & \langle b, y \rangle \\ \text{s.t.} \quad & \mathcal{A}^*y + s = c, \\ & s \in \mathcal{K}^* \end{aligned}$$

or equivalently

$$\begin{aligned} \max \quad & \langle b, y \rangle \\ \text{s.t.} \quad & c - \mathcal{A}^*y \in \mathcal{K}^* \end{aligned}$$

No one will question the above (Lagrange) dual!

Now consider the convex composite quadratic programming (CCQP)

$$\min_{x \in \mathcal{X}} \left\{ \frac{1}{2} \langle x, Qx \rangle + \langle c, x \rangle + \psi(x) \mid Ax = b \right\}$$

- $\mathcal{X}$  and  $\mathcal{Y}$  are two finite-dimensional real Hilbert spaces
- $\psi : \mathcal{X} \rightarrow (-\infty, +\infty]$  is a closed proper convex function, e.g.,  $\psi(\cdot) = \delta_P(\cdot)$ , the indicator function over a closed convex set  $P$
- $Q : \mathcal{X} \rightarrow \mathcal{X}$  satisfying  $Q = Q^*$ ,  $Q \succeq 0$
- $A : \mathcal{X} \rightarrow \mathcal{Y}$  is a given linear mapping
- $b \in \mathcal{Y}$  and  $c \in \mathcal{X}$  are given vectors

Equivalently,

$$\min_{u, x \in \mathcal{X}} \left\{ \psi(u) + \langle c, x \rangle + \frac{1}{2} \langle x, \mathcal{Q}x \rangle \mid u - x = 0, \quad \mathcal{A}x = b \right\}$$

The corresponding Lagrange function is

$$\begin{aligned} L(u, x; y, s) &:= \psi(u) + \langle c, x \rangle + \frac{1}{2} \langle x, \mathcal{Q}x \rangle + \langle s, u - x \rangle + \langle y, b - \mathcal{A}x \rangle \\ &= \langle y, b \rangle + \psi(u) + \langle s, u \rangle + \frac{1}{2} \langle x, \mathcal{Q}x \rangle + \langle x, c - \mathcal{A}^*y - s \rangle \end{aligned}$$

and the Lagrange dual of CCQP takes the form of

$$\max_{y \in \mathcal{Y}, s \in \mathcal{X}} \left\{ \inf_{u \in \mathcal{X}, x \in \mathcal{X}} L(u, x; y, s) \right\}$$

or

$$\max_{y \in \mathcal{Y}, s \in \mathcal{X}} \left\{ \langle y, b \rangle + \inf_{u \in \mathcal{X}} \left\{ \psi(u) + \langle s, u \rangle \right\} + \inf_{x \in \mathcal{X}} \left\{ \frac{1}{2} \langle x, \mathcal{Q}x \rangle + \langle x, c - \mathcal{A}^*y - s \rangle \right\} \right\}$$



By simplifications, we get the following Lagrange dual

$$\max_{y \in \mathcal{Y}, s \in \mathcal{X}} \left\{ -\psi^*(-s) + \langle y, b \rangle + \theta(y, s) \right\},$$

where

$$\theta(y, s) := \inf_{x \in \mathcal{X}} \left\{ \frac{1}{2} \langle x, Qx \rangle + \langle x, c - \mathcal{A}^*y - s \rangle \right\}$$

and  $\psi^*(\cdot)$  is the Fenchel conjugate of  $\psi$  defined by

$$\psi^*(z) := \sup_{u \in \mathcal{X}} \{ \langle z, u \rangle - \psi(u) \}.$$

Since the computation of  $\theta(y, s)$  is complicated, instead one normally considers the following **Wolfe dual** [Wolfe, Quart. Appl. Math 1961]

$$\max_{s \in \mathcal{X}, x \in \mathcal{X}, y \in \mathcal{Y}} \left\{ -\psi^*(-s) - \frac{1}{2} \langle x, Qx \rangle + \langle y, b \rangle \mid s - Qx + \mathcal{A}^*y = c \right\}.$$

Note that in the Wolfe dual (in the minimization format)

$$\min_{s \in \mathcal{X}, x \in \mathcal{X}, y \in \mathcal{Y}} \left\{ \psi^*(-s) + \frac{1}{2} \langle x, Qx \rangle - \langle y, b \rangle \mid s - Qx + A^*y = c \right\},$$

the primal variable  $x$  is also involved. But more seriously, its solution set, if nonempty, is always unbounded as long as  $Q \neq 0$  (the null space of  $Q$  is uncontrollable). It differs substantially from linear conic programming

- For Linear Semidefinite Programming:
  - The **primal** (dual, respectively) constraint nondegeneracy is equivalent to the **dual** (primal, respectively) strong second order sufficient condition (SSOSC) [Chan and S., SIOPT (2008)]
  - The **generalized Clarke Jacobian** (or the B-subdifferential) of the nonsmooth Karush- Kuhn-Tucker (KKT) solution mapping is nonsingular iff the primal and dual constraint non-degeneracies hold [Chan and S., SIOPT (2008)]
  - The **primal** (dual, respectively) strict Robinson constraint qualification is equivalent to the **dual** (primal, respectively) second order sufficient condition (SOSC) [Chao Ding, S., Liwei Zhang, SIOPT (2017)]
- The software SDPNAL for SDP [Xinyuan Zhao, S., Toh, SIOPT 2010] is applied to the dual (low rank or high rank property used)
- Convex quadratic SDP with the Wolfe dual: the above key theoretical connections hold only if  $Q \succ 0$  [Houduo Qi, MOR (2009)]
- Hard, if possible at all, to design an analogue of the software SDPNAL for the Wolfe dual if  $Q \not\succeq 0$

Our remedy is to consider the following restricted Wolfe dual (in the minimization format)

$$\min_{s \in \mathcal{X}, x' \in \mathcal{X}', y \in \mathcal{Y}} \left\{ \psi^*(-s) + \frac{1}{2} \langle x', Qx' \rangle - \langle y, b \rangle \mid s - Qx' + \mathcal{A}^*y = c \right\},$$

where  $\mathcal{X}'$  is the range space of  $Q$ , i.e.,

$$\mathcal{X}' := \text{Range}(Q).$$

One can easily check that  $Q : \mathcal{X}' \rightarrow \mathcal{X}'$  is self-adjoint and **positive definite** even if  $Q : \mathcal{X} \rightarrow \mathcal{X}$  is not positive definite. Note that if  $Q = 0$ , then  $\mathcal{X}' = \{0\}$  (in this case  $Q$  is still positive definite on  $\mathcal{X}'$  – using definition to verify it!).

Also note that  $x'$  in the dual is different from  $x$  in the primal, which does not need to stay in  $\text{Range}(Q)$ .

Different from the Wolfe dual, the **restricted Wolfe dual** possesses the following nice properties:

- It unifies with linear conic programming
- The wonderful equivalent connections between the primal and dual for SDP are now kept [Ying Cui, Chao Ding, 2019 (under revision)]
- The software QSDPNAL<sup>5</sup> for the restricted Wolfe dual of the convex quadratic SDP is developed [Xudong Li, S., Toh, MPC 2018]
- In a word, both theory and computation favor the restricted Wolfe dual

Too good to be true? — the key is to keep  $x'$  in Range (Q). But, how?  
Implementable?

---

<sup>5</sup>Publicly available at <https://blog.nus.edu.sg/mattohkc/software/qsdpnal/> and <https://github.com/MatOpt/QSDPNAL>

For given  $\sigma > 0$ , the augmented Lagrange function of the restricted Wolfe dual of the CCQP can be written as

$$\begin{aligned} L_\sigma(s, x', y; x) &:= \psi^*(-s) + \frac{1}{2} \langle x', Qx' \rangle - \langle y, b \rangle \\ &\quad + \langle x, s - Qx' + \mathcal{A}^*y - c \rangle + \frac{\sigma}{2} \|s - Qx' + \mathcal{A}^*y - c\|^2, \end{aligned}$$

which, fixing the dual variable  $x$ , is a proper closed convex (non-separable) function in the first block variable  $s$  plus a convex quadratic function in terms of  $(s, x', y)$ .

Note that  $y$  can be further split into many pieces as you please:  $y = [y^1; y^2; \dots; y^p]$ . No need to decompose  $s$  while  $x'$  must be kept as one block if the range space is used: a total of  $1 + (1 + p)$  blocks

# Augmented Lagrange function of the primal CCQP

For given  $\sigma > 0$ , the augmented Lagrange function of the CCQP (primal)

$$\min_{u, x \in \mathcal{X}} \left\{ \psi(u) + \langle c, x \rangle + \frac{1}{2} \langle x, \mathcal{Q}x \rangle \mid u - x = 0, \quad \mathcal{A}x = b \right\}$$

takes the form of

$$\begin{aligned} L_\sigma(u, x; y, s) &:= \psi(u) + \frac{1}{2} \langle x, \mathcal{Q}x \rangle + \langle c, x \rangle + \langle y, b - \mathcal{A}x \rangle + \langle s, u - x \rangle \\ &\quad + \frac{\sigma}{2} \|b - \mathcal{A}x\|^2 + \frac{\sigma}{2} \|u - x\|^2, \end{aligned}$$

which, fixing the dual variables  $y$  and  $s$ , is a proper closed convex function in the first block variable  $u$  plus a convex quadratic (non-separable) function in terms of  $(u, x)$ .

- Note that  $x$  can be further split into as many pieces as you like:  $x = [x^1; x^2; \dots; x^q]$ . There are a total of  $1 + q$  blocks.
- The introduction of  $u$  is not only necessary for the restricted Wolfe dual but also crucial for computations: not a computationally good idea to split  $x$  into multiple parts without  $u$  even if  $\psi(\cdot)$  is separable.

- Note that the augmented Lagrange function for the CCQP in the primal form does not contain a nonsmooth term for  $x$  even if it is non-separable for  $x$ .
- The (mysterious?) forms of the two augmented Lagrange functions lead to the discovery of the symmetric Gauss-Seidel decomposition theorem!



Actually, we can consider more general convex composite quadratic programming (CCQP)

$$\min_{x \in \mathcal{X}} \left\{ \psi(x) + \frac{1}{2} \langle x, \mathcal{Q}x \rangle - \langle c, x \rangle \mid \mathcal{A}_E x = b_E, \mathcal{A}_I x - b_I \in \mathcal{K} \right\}$$

- $\psi : \mathcal{X} \rightarrow (-\infty, +\infty]$  is a closed proper convex function [simple]
- $\mathcal{Q} : \mathcal{X} \rightarrow \mathcal{X}$  satisfying  $\mathcal{Q} = \mathcal{Q}^*$ ,  $\mathcal{Q} \succeq 0$
- $\mathcal{A}_E : \mathcal{X} \rightarrow \mathcal{Z}_1$  and  $\mathcal{A}_I : \mathcal{X} \rightarrow \mathcal{Z}_2$ , given linear mappings
- $b = (b_E; b_I) \in \mathcal{Z} := \mathcal{Z}_1 \times \mathcal{Z}_2$ , given vector
- $c \in \mathcal{X}$  is given.
- $\mathcal{K} \subseteq \mathcal{Z}_2$  is a closed convex set (cone) [simple]
- $\mathcal{X}$ ,  $\mathcal{Z}_1$ , and  $\mathcal{Z}_2$  are finite-dimensional real Hilbert spaces

Equivalently,

$$\min_{x \in \mathcal{X}, x' \in \mathcal{Z}_2} \left\{ \psi(x) + \delta_{\mathcal{K}}(x') + \frac{1}{2} \langle x, \mathcal{Q}x \rangle - \langle c, x \rangle \mid \begin{pmatrix} \mathcal{A}_E & 0 \\ \mathcal{A}_I & -\mathcal{I} \end{pmatrix} \begin{pmatrix} x \\ x' \end{pmatrix} = b \right\},$$

whose **restricted Wolfe dual** (in the minimization format) is

$$\min_{\substack{s \in \mathcal{Y}, z \in \mathcal{Z} \\ y' \in \text{Range}(\mathcal{Q})}} \left\{ p(s) + \frac{1}{2} \langle y', \mathcal{Q}y' \rangle - \langle b, z \rangle \mid s + \begin{pmatrix} \mathcal{Q} \\ 0 \end{pmatrix} y' - \begin{pmatrix} \mathcal{A}_E^* & \mathcal{A}_I^* \\ 0 & -\mathcal{I} \end{pmatrix} z = \begin{pmatrix} c \\ 0 \end{pmatrix} \right\}$$

- $s := (u, v) \in \mathcal{Y} := \mathcal{X} \times \mathcal{Z}_2$
- $p(s) := p(u, v) = \psi^*(u) + \delta_{\mathcal{K}}^*(v)$
- $\delta_{\mathcal{K}}(\cdot)$  is the indicator function over  $\mathcal{K}$

# A general form of symmetric Gauss-Seidel iteration

Consider the **block** vector

$\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_s) \in \mathcal{X} := \mathcal{X}_1 \times \mathcal{X}_2 \times \dots \times \mathcal{X}_s$ . Given a positive semidefinite linear operator  $\mathcal{Q}$  such that

$$\mathcal{Q}\mathbf{x} \equiv \begin{pmatrix} \mathcal{Q}_{11} & \mathcal{Q}_{12} & \cdots & \mathcal{Q}_{1s} \\ \mathcal{Q}_{12}^* & \mathcal{Q}_{22} & \cdots & \mathcal{Q}_{2s} \\ \vdots & \vdots & \ddots & \vdots \\ \mathcal{Q}_{1s}^* & \mathcal{Q}_{2s}^* & \cdots & \mathcal{Q}_{ss} \end{pmatrix} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_s \end{pmatrix}, \quad \mathcal{Q}_{ii} \succ \mathbf{0}.$$

Let  $p : \mathcal{X}_1 \rightarrow (-\infty, +\infty]$  be a given closed proper convex function. Let the quadratic function

$$q(\mathbf{x}) := \frac{1}{2} \langle \mathbf{x}, \mathcal{Q}\mathbf{x} \rangle - \langle \mathbf{r}, \mathbf{x} \rangle.$$

Consider the following **block decomposition**:

$$U\mathbf{x} \equiv \begin{pmatrix} 0 & Q_{12} & \cdots & Q_{1s} \\ & \ddots & & \vdots \\ & & \ddots & Q_{(s-1)s} \\ & & & 0 \end{pmatrix} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_s \end{pmatrix}.$$

Then  $Q = U^* + D + U$ , where  $D\mathbf{x} = (Q_{11}\mathbf{x}_1, \dots, Q_{ss}\mathbf{x}_s)$ .

Let  $\hat{\delta} \equiv (\hat{\delta}_1, \dots, \hat{\delta}_s)$  and  $\delta^+ \equiv (\delta_1^+, \dots, \delta_s^+)$  with  $\hat{\delta}_1 = \delta_1^+$  being given error tolerance vectors. Define

$$\Delta(\hat{\delta}, \delta^+) := \delta^+ + UD^{-1}(\delta^+ - \hat{\delta}), \quad \mathcal{T} := UD^{-1}U^* \text{ (sGS decomp. op.)}.$$

Note that  $\mathcal{T} \succeq 0$  is NOT positive definite. Let  $\mathbf{y} \in \mathcal{X}$  be given. Define

$$\mathbf{x}^+ := \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ p(\mathbf{x}_1) + q(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_{\mathcal{T}}^2 - \langle \Delta(\hat{\delta}, \delta^+), \mathbf{x} \rangle \right\}. \quad (1)$$

(1) looks complicated, but is much easier to solve!

Theorem (Xudong Li, S., Toh, MP 2019)

Given  $\mathbf{y}$ . For  $i = s, \dots, 2$ , define

$$\begin{aligned}\hat{\mathbf{x}}_i &:= \arg \min_{\mathbf{x}_i} \{ p(\mathbf{y}_1) + q(\mathbf{y}_{\leq i-1}, \mathbf{x}_i, \hat{\mathbf{x}}_{\geq i+1}) - \langle \hat{\delta}_i, \mathbf{x}_i \rangle \} \\ &= \mathcal{Q}_{ii}^{-1} (\mathbf{r}_i + \hat{\delta}_i - \sum_{j=1}^{i-1} \mathcal{Q}_{ji}^* \mathbf{y}_j - \sum_{j=i+1}^s \mathcal{Q}_{ij} \hat{\mathbf{x}}_j)\end{aligned}$$

computed in the *backward GS cycle*. **The optimal solution  $\mathbf{x}^+$  in (1) can be obtained exactly via**

$$\begin{aligned}\mathbf{x}_1^+ &= \arg \min_{\mathbf{x}_1} \{ p(\mathbf{x}_1) + q(\mathbf{x}_1, \hat{\mathbf{x}}_{\geq 2}) - \langle \delta_1^+, \mathbf{x}_1 \rangle \}, \\ \mathbf{x}_i^+ &= \arg \min_{\mathbf{x}_i} \{ p(\mathbf{x}_1^+) + q(\mathbf{x}_{\leq i-1}^+, \mathbf{x}_i, \hat{\mathbf{x}}_{\geq i+1}) - \langle \delta_i^+, \mathbf{x}_i \rangle \} \\ &= \mathcal{Q}_{ii}^{-1} (\mathbf{r}_i + \delta_i^+ - \sum_{j=1}^{i-1} \mathcal{Q}_{ji}^* \mathbf{x}_j^+ - \sum_{j=i+1}^s \mathcal{Q}_{ij} \hat{\mathbf{x}}_j), \quad i \geq 2,\end{aligned}$$

where  $\mathbf{x}_i^+$ ,  $i = 1, 2, \dots, s$ , is computed in the *forward GS cycle*.

Reduces to the classical block sGS if both  $p(\cdot) \equiv 0$  and  $\delta = 0$ .

**Caution:** this theorem (symmetric property) does not hold for GS even  $s = 2$ .

Consider the convex optimization model:

$$\begin{aligned} \min \quad & \theta(y_1) + f(y_1, y_2, \dots, y_s) \\ \text{s.t.} \quad & \mathcal{A}_1^* y_1 + \mathcal{A}_2^* y_2 + \dots + \mathcal{A}_s^* y_s = c. \end{aligned} \tag{2}$$

Linear mappings:  $\mathcal{A}_i, i = 1, \dots, s, \mathcal{A}^* y = \sum_{i=1}^s \mathcal{A}_i^* y_i, y := (y_1, \dots, y_s)$ .  
 Closed proper convex function  $\theta : \mathcal{Y}_1 \rightarrow (-\infty, +\infty]$  and convex quadratic function  $f(y) = \frac{1}{2} \langle y, \mathcal{Q}y \rangle - \langle b, y \rangle$ . Then, (3) can be written compactly as

$$\min \{ \theta(y_1) + f(y) \mid \mathcal{A}^* y = c \},$$

which is a very general CCQP.

Given  $\sigma > 0$ , the augmented Lagrangian function of the CCQP is

$$\mathcal{L}_\sigma(y; x) = \theta(y_1) + f(y) + \underbrace{\langle x, \mathcal{A}^* y - c \rangle + \frac{\sigma}{2} \|\mathcal{A}^* y - c\|^2}_{\text{quadratic}}.$$

The proximal augmented Lagrangian method (pALM) for the CCQP:

---

Given  $(y^0, x^0)$  in the domain and  $\tau \in (0, 2)$ . For  $k = 0, 1, \dots$

**Step 1.**  $y^{k+1} \approx \arg \min \mathcal{L}_\sigma(y; x^k) + \frac{1}{2} \|y - y^k\|_{\mathcal{T}}^2$

$$= \arg \min_y \left\{ \theta(y_1) + f(y) + \langle x^k, \mathcal{A}^*y - c \rangle + \frac{\sigma}{2} \|\mathcal{A}^*y - c\|^2 + \frac{1}{2} \|y - y^k\|_{\mathcal{T}}^2 \right\}.$$

**Step 2.**  $x^{k+1} = x^k + \tau \sigma (\mathcal{A}^*y^{k+1} - c)$ .

---

- $\mathcal{T}$  is the block sGS decomposition operator of  $Q + \sigma \mathcal{A} \mathcal{A}^*$ , which does not need to be formulated explicitly. Note that  $\mathcal{T} \succeq 0$  but  $\mathcal{T} \neq 0$ . So it is not a classical pALM, but a “semiproximal” ALM.
- $y^{k+1}$  is obtained via the inexact block sGS procedure [ $s$  blocks in total].
- In practice, the dual step-length  $\tau$  is often chosen in [1.618, 1.95], e.g.,  $\tau = 1.9$ .

Consider the LCCQP with **two non-quadratic blocks**:

$$\begin{aligned} \min \quad & \theta(y_1) + q(y_1, \dots, y_s) + \phi(z_1) + h(z_1, \dots, z_t) \\ \text{s.t.} \quad & \mathcal{A}_1^* y_1 + \dots + \mathcal{A}_s^* y_s + \mathcal{B}_1^* z_1 + \dots + \mathcal{B}_t^* z_t = c \end{aligned} \quad (3)$$

linear mappings:  $\mathcal{A}_i, i = 1, \dots, s, \mathcal{A}^* y = \sum_{i=1}^s \mathcal{A}_i^* y_i, y := (y_1, \dots, y_s)$

linear mappings:  $\mathcal{B}_j, j = 1, \dots, t, \mathcal{B}^* z = \sum_{i=1}^t \mathcal{B}_i^* z_i, z := (z_1, \dots, z_t)$

closed proper convex functions  $\theta : \mathcal{Y}_1 \rightarrow (-\infty, \infty], \phi : \mathcal{Z}_1 \rightarrow (-\infty, \infty]$

convex quadratics:  $q(y) = \frac{1}{2} \langle y, Qy \rangle - \langle b, y \rangle, h(z) = \frac{1}{2} \langle z, Hz \rangle - \langle d, z \rangle$

Write (3) compactly as

$$\min \{ \theta(y_1) + q(y) + \phi(z_1) + h(z) \mid \mathcal{A}^* y + \mathcal{B}^* z = c \}.$$

Given  $\sigma > 0$ , the associate augmented Lagrangian function is

$$\mathcal{L}_\sigma(y, z; x) = \theta(y_1) + q(y) + \phi(z_1) + h(z) + \frac{\sigma}{2} \|\mathcal{A}^* y + \mathcal{B}^* z - c + x/\sigma\|^2 - \frac{\|x\|^2}{2\sigma}$$



Given  $(y^0, z^0, x^0)$  in the domain. Iterate

$$\text{Step 1. } y^{k+1} \approx \arg \min_y \left\{ \mathcal{L}_\sigma(y, z^k; x^k) + \frac{1}{2} \|y - y^k\|_{\mathcal{T}}^2 \right\}$$

$$= \arg \min_{y=(y_1, y_2, \dots, y_s)} \left\{ \theta(y_1) + \frac{1}{2} \langle y, (\mathcal{Q} + \sigma \mathcal{A} \mathcal{A}^*) y \rangle - \langle r^k, y \rangle + \frac{1}{2} \|y - y^k\|_{\mathcal{T}}^2 \right\}.$$

$$\text{Step 2. } z^{k+1} \approx \arg \min_z \left\{ \mathcal{L}_\sigma(y^{k+1}, z; x^k) + \frac{1}{2} \|z - z^k\|_{\mathcal{S}}^2 \right\}$$

$$= \arg \min_{z=(z_1, z_2, \dots, z_t)} \left\{ \phi(z_1) + \frac{1}{2} \langle z, (\mathcal{H} + \sigma \mathcal{B} \mathcal{B}^*) z \rangle - \langle s^k, z \rangle + \frac{1}{2} \|z - z^k\|_{\mathcal{S}}^2 \right\}.$$

**Step 3.**  $x^{k+1} = x^k + \tau \sigma (\mathcal{A}^* y^{k+1} + \mathcal{B}^* z^{k+1} - c)$ , where step-length  $\tau \in (0, \frac{1+\sqrt{5}}{2})$ .

- $\mathcal{T}$  = block sGS operator of  $\mathcal{Q} + \sigma \mathcal{A} \mathcal{A}^*$ ,  
 $\mathcal{S}$  = block sGS operator of  $\mathcal{H} + \sigma \mathcal{B} \mathcal{B}^*$ .
- $y^{k+1} = (y_1, y_2, \dots, y_s)^{k+1}$  is obtained via one cycle of the inexact block sGS iteration. Similarly for  $z^{k+1} = (z_1, z_2, \dots, z_t)^{k+1}$ .
- sGS-pADMM has well developed convergence guarantee.

Define  $\mathcal{N} = \{X \in \mathbb{S}^n \mid X \geq 0\}$  (cone of nonnegative matrices)

$$\min \left\{ \langle C, X \rangle \mid \mathcal{A}(X) = b, X \in \mathbb{S}_+^n, X \in \mathcal{N} \right\}$$

$$\text{(dual)} \quad - \min \left\{ \delta_{\mathbb{S}_+^n}(S) - \langle b, y \rangle + \delta_{\mathcal{N}}(Z) \mid S + \mathcal{A}^*y + Z = C \right\}$$

Input  $(y_0, S_0; X_0)$ . For  $l = 0, 1, \dots$ , let  $\widehat{C}_l = C - \sigma^{-1}X_l$

$$\text{[1a]} \quad \widehat{y}_{l+1} \approx \operatorname{argmin}_{y \in \mathbb{R}^m} \{ \mathcal{L}_\sigma(y, S_l, Z_l; X_l) \} \rightarrow \text{solve } \mathcal{A}\mathcal{A}^*y = \text{rhs}$$

$$\text{[1b]} \quad S_{l+1} = \operatorname{argmin}_{S \in \mathbb{S}_+^n} \{ \mathcal{L}_\sigma(\widehat{y}_{l+1}, S, Z_l; X_l) \} = \Pi_{\mathbb{S}_+^n}(\widehat{C}_l - \mathcal{A}^*\widehat{y}_{l+1} - Z_l)$$

$$\text{[1c]} \quad y_{l+1} \approx \operatorname{argmin}_{y \in \mathbb{R}^m} \{ \mathcal{L}_\sigma(y, S_{l+1}, Z_l; X_l) \} \rightarrow \text{solve } \mathcal{A}\mathcal{A}^*y = \text{rhs}$$

$$\text{[2]} \quad Z_{l+1} = \operatorname{argmin}_{Z \in \mathcal{N}} \{ \mathcal{L}_\sigma(y_{l+1}, S_{l+1}, Z; X_l) \} = \Pi_{\mathcal{N}}(\widehat{C}_l - \mathcal{A}^*y_{l+1} - S_{l+1})$$

$$\text{[3]} \quad X_{l+1} = X_l + \tau\sigma(\mathcal{A}^*y_{l+1} + S_{l+1} + Z_{l+1} - C), \quad \tau \in (0, \frac{1+\sqrt{5}}{2})$$

sGS-pADMM is a convergent enhancement [S., Toh, Liuqin Yang, SIOPT 2015] of the sequential Gauss-Seidel ADMM whose convergence is not guaranteed. [one can swop the positions of  $S$  and  $Z$ ]

Number of problems solved to the accuracy of  $10^{-6}$  in relative KKT residual.

problem set (No.) \ solver	sGS-ADMM	SDPAD	2EBD	ADMM3g
$\theta_+$ (58)	58	58	56	54
FAP ( 7)	7	7	7	7
QAP (95)	39	30	16	28
BIQ (134)	134	134	134	130
RCP (120)	120	114	109	113
Total (414)	358	343	322	332

## Comparison of sGS-ADMM, SDPAD, ADMM3g and 2EBD

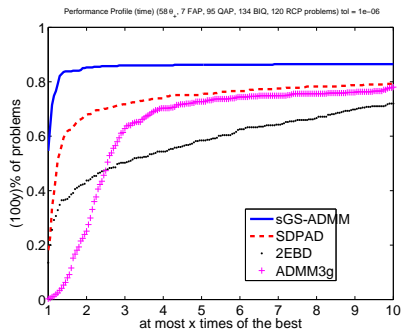


Figure: Performance profiles (time).

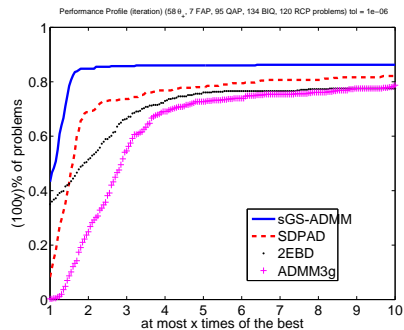


Figure: Performance profiles (iteration).

- In the sGS procedure, the coefficient matrices of the linear systems to be solved only need to be **factorized once at the start of the procedure**. The additional costs of the repetitions are minimal and can be offset by the larger step length  $\tau \in (0, 2)$  or  $\tau \in (0, \frac{1+\sqrt{5}}{2})$ .
- There are many applications that can be “solved” via block **sGS + pALM** or block **sGS+ pADMM** if the solution accuracy is not a big concern.
- More extensions can be done. For example, for convex quadratic semidefinite programming, and other convex composite conic programming problems.
- **To make the algorithms even faster, we often introduce indefinite proximal terms with guaranteed convergence [Liang Chen, S., Toh, Ning Zhang, JCM 2019]**
- **We need something more than the “sGS + pALM” (or “pADMM)”** – semismooth Newton methods (SDPNAL+ solves all the above 414 examples)

# One very simple example: a convex QP with Birkhoff

Convex QP:

$$(\mathbf{P}) \quad \min \left\{ \frac{1}{2} \langle X, \mathcal{Q}X \rangle + \langle G, X \rangle \mid X \in \mathfrak{B}_n \right\},$$

self-adjoint linear operator  $\mathcal{Q} \succeq 0$  and the Birkhoff polytope:

$$\mathfrak{B}_n := \{X \in \mathbb{R}^{n \times n} \mid Xe = e, X^T e = e, X \geq 0\}$$

$e \in \mathbb{R}^n$ : the vector of all ones.

$$(\mathbf{D}) \quad \min \left\{ \delta_{\mathfrak{B}_n}^*(Z) + \frac{1}{2} \langle W, \mathcal{Q}W \rangle \mid Z + \mathcal{Q}W + G = 0, W \in \text{Range}(\mathcal{Q}) \right\}$$

$\delta_{\mathfrak{B}_n}^*$ : the conjugate of the indicator function  $\delta_{\mathfrak{B}_n}$  [Li, S., Toh, MP 2020]

ALM function for  $(\mathbf{D})$ , given  $\sigma > 0$

$$\begin{aligned} \mathcal{L}_\sigma(Z, W; X) &= \delta_{\mathfrak{B}_n}^*(Z) + \frac{1}{2} \langle W, \mathcal{Q}W \rangle - \langle X, Z + \mathcal{Q}W + G \rangle \\ &\quad + \frac{\sigma}{2} \|Z + \mathcal{Q}W + G\|^2 \end{aligned}$$

**Algorithm ALM: An augmented Lagrangian method for (D).**

Given  $\sigma_0 > 0$ , iterates  $k = 0, 1, \dots$

Step 1. Compute

$$(Z^{k+1}, W^{k+1}) \approx \operatorname{argmin} \left\{ \begin{array}{l} \Psi_k(Z, W) := \mathcal{L}_{\sigma_k}(Z, W; X^k) \\ | (Z, W) \in \mathbb{R}^{n \times n} \times \operatorname{Range}(\mathcal{Q}) \end{array} \right\}.$$

Step 2. Compute

$$X^{k+1} = X^k - \sigma_k(Z^{k+1} + \mathcal{Q}W^{k+1} + G).$$

Update  $\sigma_{k+1} \uparrow \sigma_\infty \leq \infty$ .

Convex **piecewise linear-quadratic** minimization (SDP is more complicated):

**error bound holds**  $\implies$  ALM converges asymptotically **superlinearly**

# Semismooth Newton-CG method for inner problem

For any  $W \in \text{Range}(\mathcal{Q})$ ,

$$\psi(W) := \inf_Z \mathcal{L}_\sigma(Z, W; \hat{X}), \quad Z(W) := \hat{X} - \sigma(\mathcal{Q}W + G)$$

Subproblem solution  $(\bar{Z}, \bar{W})$ :

$$\begin{aligned}\bar{W} &= \arg \min \{ \psi(W) \mid W \in \text{Range}(\mathcal{Q}) \}, \\ \bar{Z} &= \sigma^{-1}(Z(\bar{W}) - \Pi_{\mathfrak{B}_n}(Z(\bar{W})))\end{aligned}$$

For all  $W \in \text{Range}(\mathcal{Q})$ ,

$$\nabla \psi(W) = \mathcal{Q}(W - \Pi_{\mathfrak{B}_n}(Z(W)))$$

Semismooth Newton CG solves nonsmooth piecewise affine equation

$$\nabla \psi(W) = 0, \quad W \in \text{Range}(\mathcal{Q}).$$



Given  $\widehat{W}$ , linear operator  $\mathcal{M} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n}$

$$\mathcal{M}(\Delta W) := \mathcal{Q}(\mathcal{I} + \sigma P_{HS} \mathcal{Q}) \Delta W, \quad \forall \Delta W \in \mathbb{R}^{n \times n}$$

$P_{HS}$ : the HS-Jacobian of  $\Pi_{\mathfrak{B}_n}$  at  $Z(\widehat{W})$

$j$ -th iter., solve linear system (CG)

$$\mathcal{M}_j dW + \nabla \psi(W^j) = 0, \quad dW \in \text{Range}(\mathcal{Q})$$

Global convergence: Line search (using  $\psi(W)$ )

Local convergence:

positive definiteness of  $\mathcal{M}$  on  $\text{Range}(\mathcal{Q}) \implies$  at least **superlinear**

Given  $A, B \in \mathcal{S}^n$ , quadratic assignment problem (QAP):

$$\min\{\langle X, AXB \rangle \mid X \in \{0, 1\}^{n \times n} \cap \mathfrak{B}_n\}$$

Convex relaxation [Anstreicher et al. MP, 2001]:

$$\min\{\langle X, QX \rangle \mid X \in \mathfrak{B}_n\}$$

Self-adjoint linear operator  $Q(X) := AXB - SX - XT$ ,  $Q \succeq 0$

Matrices  $S, T \in \mathcal{S}^n$  obtained from [Anstreicher et al. MP, 2001]

Relative KKT residual:

$$\eta = \frac{\|X - \Pi_{\mathfrak{B}_n}(X - QX)\|}{1 + \|X\| + \|QX\|}$$

Matrices  $A, B$  from QAPLIB

# Numerical results for QAP

“a”: Gurobi, “b”: ALM

		iter		$\eta$	time
problem	$n$	a	b (itersub)	a b	a b
lipa80a	80	11	25 (68)	1.3-6   7.3-8	2:46   01
lipa90a	90	11	20 (54)	2.7-6   8.8-8	5:32   01
sko100a	100	14	26 (95)	8.5-6   8.5-8	2:06   11
tai100a	100	11	18 (52)	1.3-6   9.5-8	10:31   02
tai100b	100	11	27 (98)	1.3-6   9.1-8	10:31   13
tai80b	80	11	27 (98)	1.2-6   8.5-8	2:36   07
tai256c	256	*	2 ( 4)	*   2.1-16	*   00
tai150b	150	19	27 (94)	4.3-7   9.3-8	2:46:17   13
tho150	150	16	24 (96)	5.6-6   9.9-8	18:52   22

“\*”: Gurobi out of memory (128 G RAM)

“tai150b”: Gurobi reports error, “small positive term” needed

- SDPNAL/SDPNAL+ and QSDPNAL: two-phase augmented Lagrangian methods for solving large scale SDP and convex quadratic SDPs, which are [publically available](#).
- Primal constraint non-degeneracy is the key for using the semismooth Newton methods successfully (e.g., the NCM problem); the quadratic growth condition is crucial for the fast convergence of the (dual) ALM; and the primal-dual errors bounds can easily fail<sup>6</sup>
- We are still at the very early stages of solving large scale SDPs and beyond: more theory, more algorithms, and more software packages are needed.

---

<sup>6</sup>Ying Cui and S., and Toh, “On the R-superlinear convergence of the KKT residuals generated by the augmented Lagrangian method for convex composite conic programming, *Mathematical Programming* 178 (2019) 381–415

**Thank you for your attention!**