

**Convex Optimization Methods for
Dimension Reduction and
Coefficient Estimation in
Multivariate Linear Regression**

**Renato D.C. Monteiro
Georgia Tech**

FODAVA REVIEW MEETING

Atlanta, USA

Dec. 3, 2009

MOTIVATION

Compressed sensing: Candes et al. (2006) have shown that a sparse signal can be recovered by solving a non-smooth optimization problem of the form

$$\min\{\|\mathbf{x}\|_0 : \mathbf{Ax} = \mathbf{b}\} \quad (*)$$

where \mathbf{A} is $l \times p$ -matrix, $\mathbf{x} \in \mathbb{R}^p$ and $\mathbf{b} \in \mathbb{R}^l$. Here $\|\mathbf{x}\|_0$ denotes the number of nonzero components of \mathbf{x} .

Under some conditions on \mathbf{A} , they have shown that $(*)$ is also equivalent to the convex program

$$\min\{\|\mathbf{x}\|_1 : \mathbf{Ax} = \mathbf{b}\} \quad (**)$$

where $\|\mathbf{x}\|_1 := \sum_{i=1}^p |\mathbf{x}_i|$ is the **1**-norm of \mathbf{x} .

There has been a lot of research in the optimization community to develop methods that can solve $(**)$ efficiently.

We are interested in the extensions of these problems, where the variable is now a $p \times q$ matrix.

Consider the problem

$$\min_{\mathbf{X} \in \mathbb{R}^{\mathbf{p} \times \mathbf{q}}} \{\text{rank}(\mathbf{X}) : \mathcal{A}(\mathbf{X}) = \mathbf{b}\}, \quad (*)$$

where $\mathcal{A} : \mathbb{R}^{\mathbf{p} \times \mathbf{q}} \rightarrow \mathbb{R}^{\mathbf{l}}$ is a linear map and $\mathbf{b} \in \mathbb{R}^{\mathbf{l}}$.

This problem (and its variations) has many applications (e.g., matrix completion problems, netflix problem, dimension reduction in statistics and etc.)

When \mathbf{X} is restricted to a diagonal (square) matrix, i.e., $\mathbf{X} = \text{Diag}(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^{\mathbf{p}}$, then $(*)$ reduces to

$$\min\{\|\mathbf{x}\|_0 : \mathbf{A}\mathbf{x} = \mathbf{b}\}$$

A convex approximation of $(*)$ is

$$\min_{\mathbf{X} \in \mathbb{R}^{\mathbf{p} \times \mathbf{q}}} \{\|\mathbf{X}\|_* : \mathcal{A}(\mathbf{X}) = \mathbf{b}\} \quad (**)$$

where $\|\mathbf{X}\|_*$ denotes the nuclear norm of \mathbf{X} :

$$\|\mathbf{X}\|_* := \text{Trace}[(\mathbf{X}\mathbf{X}^{\mathbf{T}})^{1/2}] = \sum_{i=1}^{\min\{\mathbf{p}, \mathbf{q}\}} \sigma_i(\mathbf{X})$$

Under suitable conditions on \mathcal{A} , Recht et al. (2007) have shown that $(*)$ and $(**)$ are equivalent.

A relaxation of (**) is the problem

$$\min_{\mathbf{X} \in \mathbb{R}^{p \times q}} \frac{1}{2} \|\mathcal{A}(\mathbf{X}) - \mathbf{b}\|_{\mathbb{F}}^2 + \lambda \|\mathbf{X}\|_*$$

where $\|\mathbf{X}\|_{\mathbb{F}}^2 := \sum_{i,j} \mathbf{X}_{ij}^2$.

Our problem of interest is the following special case of the above problem:

$$\min_{\mathbf{X} \in \mathbb{R}^{p \times q}} \frac{1}{2} \|\mathbf{A}\mathbf{X} - \mathbf{B}\|_{\mathbb{F}}^2 + \lambda \|\mathbf{X}\|_*$$

where $\mathbf{A} \in \mathbb{R}^{n \times p}$, $\mathbf{B} \in \mathbb{R}^{n \times q}$ and \mathbf{A} has l.i. columns ($n \gg p$).

This problem arises in statistics in the context of dimension reduction and coefficient estimate in multivariate linear regression.

DIMENSION REDUCTION IN STATISTICS

Assume that $\mathbf{A} \in \mathbb{R}^{n \times p}$ consists of n observations on p explanatory variables $\mathbf{a} = (\mathbf{a}_1, \dots, \mathbf{a}_p)'$ and $\mathbf{B} \in \mathbb{R}^{n \times q}$ collects the corresponding n observations on q responses $\mathbf{b} = (\mathbf{b}_1, \dots, \mathbf{b}_q)'$. Consider the multivariate linear model

$$\mathbf{B} = \mathbf{A}\mathbf{X} + \mathbf{E},$$

where $\mathbf{X} \in \mathbb{R}^{p \times q}$ is a coefficient matrix, $\mathbf{E} = (\mathbf{e}^1, \dots, \mathbf{e}^n)'$ is the regression noise, and all \mathbf{e}^i 's are independent samples of $\mathcal{N}(\mathbf{0}, \Sigma)$.

To estimate \mathbf{X} and accomplish dimension reduction, Yuan et al. proposed to solve

$$\min_{\mathbf{X}} \frac{1}{2} \|\mathbf{A}\mathbf{X} - \mathbf{B}\|_{\mathbf{F}}^2 + \lambda \|\mathbf{X}\|_* \quad (1)$$

for different $\lambda > 0$ values. The larger the scalar $\lambda > 0$, the more dimension reduction is accomplished.

Reformulations:

- cone program (includes LP, SDP)
- saddle point (min-max convex-concave) problems

Cone program (CP): Given a closed convex cone $\mathcal{K} \subseteq \mathbb{R}^n$, the CP problem is:

$$\min\{\langle \mathbf{c}, \mathbf{x} \rangle : \mathcal{A}(\mathbf{x}) = \mathbf{b}, \mathbf{x} \in \mathcal{K}\}$$

where $\mathbf{c} \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^m$ and $\mathcal{A} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a linear map. Its dual is

$$\max\{\langle \mathbf{b}, \mathbf{y} \rangle : \mathbf{c} - \mathcal{A}^*(\mathbf{y}) \in \mathcal{K}^*\}$$

where $\mathcal{K}^* := \{\mathbf{s} \in \mathbb{R}^n : \langle \mathbf{s}, \mathbf{x} \rangle \geq \mathbf{0}, \forall \mathbf{x} \in \mathcal{K}\}$.

Remark: Can be solved by interior-point (second-order) methods or by first-order methods.

SADDLE POINT OR MIN-MAX PROBLEMS

Their general form is

$$\min_{\mathbf{x} \in \mathbf{X}} \left(\mathbf{f}(\mathbf{x}) := \max_{\mathbf{y} \in \mathbf{Y}} \phi(\mathbf{x}, \mathbf{y}) \right)$$

where \mathbf{X}, \mathbf{Y} are simple closed convex sets, ϕ is convex in \mathbf{x} and concave in \mathbf{y} .

Under the assumption that $\nabla \phi$ is Lipschitz continuous, first-order methods with known iteration-complexity bounds have been developed to solve these problems:

- Nesterov's smooth or non-smooth methods and their variants;
- Korpelevich algorithm or Nemirovski's prox-mirror method

CONE PROGRAMMING REFORMULATION

Problem (1) can be reformulated as a CP problem as follows. Clearly, (1) is equivalent to

$$\min_{\mathbf{X}, t} \left\{ \frac{1}{2} \|\mathbf{AX} - \mathbf{B}\|^2 + \lambda t : \|\mathbf{X}\|_* \leq t \right\}$$

Write $\mathbf{V} \succeq \mathbf{0}$ if \mathbf{V} is symmetric and positive semidefinite. Also, let \mathcal{S}^l denote the space of $l \times l$ symm. matrices.

Proposition: Let $\mathbf{X} \in \mathbb{R}^{p \times q}$ and set $\mathbf{k} := \min\{p, q\}$ and $\mathbf{l} := p + q$. For $t \in \mathbb{R}$, we have

$$\|\mathbf{X}\|_* \leq t \Leftrightarrow \begin{cases} t - \mathbf{k}s - \text{Trace}(\mathbf{V}) & \geq 0, \\ \mathbf{V} - \mathcal{G}(\mathbf{X}) + s\mathbf{I} & \succeq 0, \\ \mathbf{V} & \succeq 0, \end{cases}$$

for some $\mathbf{V} \in \mathcal{S}^l$ and $s \in \mathbb{R}$, where $\mathcal{G} : \mathbb{R}^{p \times q} \rightarrow \mathbb{R}^{l \times l}$ is defined as

$$\mathcal{G}(\mathbf{X}) := \begin{pmatrix} \mathbf{0} & \mathbf{X}^T \\ \mathbf{X} & \mathbf{0} \end{pmatrix}$$

SADDLE-POINT REFORMULATIONS

Using the identity

$$\|\mathbf{X}\|_* = \mathbf{k} \max_{\mathbf{W} \in \Omega} \langle \mathcal{G}(\mathbf{X}), \mathbf{W} \rangle$$

where $\mathbf{k} := \min\{\mathbf{p}, \mathbf{q}\}$ and

$$\Omega := \{\mathbf{W} \in \mathcal{S}^{\mathbf{p}+\mathbf{q}} : \mathbf{0} \preceq \mathbf{W} \preceq \mathbf{I}/\mathbf{k}, \text{Trace}(\mathbf{W}) = \mathbf{1}\}$$

problem (1) can be reformulated as

$$\min_{\|\mathbf{X}\|_{\mathbf{F}} \leq \mathbf{r}} \mathbf{f}_{\mathbf{p}}(\mathbf{X}) := \max_{\mathbf{W} \in \Omega} \left\{ \frac{1}{2} \|\mathbf{A}\mathbf{X} - \mathbf{B}\|_{\mathbf{F}}^2 + \lambda \mathbf{k} \langle \mathcal{G}(\mathbf{X}), \mathbf{W} \rangle \right\},$$

where \mathbf{r} is an appropriate scalar. (Disadvantage: $\mathbf{f}_{\mathbf{p}}(\cdot)$ is non-smooth)

Instead, we consider the dual of the above problem, namely:

$$\max_{\mathbf{W} \in \Omega} \mathbf{f}_{\mathbf{d}}(\mathbf{W}) := \min_{\|\mathbf{X}\|_{\mathbf{F}} \leq \mathbf{r}} \left\{ \frac{1}{2} \|\mathbf{A}\mathbf{X} - \mathbf{B}\|_{\mathbf{F}}^2 + \lambda \mathbf{k} \langle \mathcal{G}(\mathbf{X}), \mathbf{W} \rangle \right\},$$

(Advantage: $\mathbf{f}_{\mathbf{d}}(\mathbf{W})$ has Lipschitz continuous gradient.)

We then apply a variant of Nesterov's smooth method to solve the latter MAX-MIN reformulation of (1).

Proposition: Given $\epsilon > 0$, Nesterov's smooth method finds an ϵ -optimal solution of the MAX-MIN formulation in a number of iterations which does not exceed

$$\frac{2\lambda\|(\mathbf{A}^T\mathbf{A})^{-1/2}\|}{\sqrt{\epsilon}}\sqrt{\mathbf{k}\log\left(\frac{\mathbf{p}+\mathbf{q}}{\mathbf{k}}\right)},$$

where $\mathbf{k} := \min\{\mathbf{p}, \mathbf{q}\}$.

Note: The complexity of solving the corresponding dual MIN-MAX reformulation of (1) is $\mathcal{O}(1/\epsilon)$ instead of $\mathcal{O}(1/\sqrt{\epsilon})$ as above.

COMPUTATIONAL RESULTS

The entries of $\mathbf{A} \in \mathfrak{R}^{n \times p}$ and $\mathbf{B} \in \mathfrak{R}^{n \times q}$, with $p = 2q$ and $n = 10q$, were uniformly generated in $[0, 1]$. The accuracy in the table below is $\epsilon = 10^{-1}$.

Problem (p, q)	# of Iterations		CPU Time	
	MIN-MAX	MAX-MIN	MIN-MAX	MAX-MIN
(200, 100)	610	1	29.60	0.91
(400, 200)	1310	1	432.92	8.36
(600, 300)	2061	1	2155.76	31.23
(800, 400)	2848	1	7831.09	76.75
(1000, 500)	3628	1	21128.70	156.68
(1200, 600)	4436	1	47356.32	276.64
(1400, 700)	5280	1	98573.73	456.61
(1600, 800)	6108	1	176557.49	699.47

COMPUTATIONAL RESULTS

The tables below compare the MAX-MIN formulation with the cone programming reformulation. The accuracy is $\epsilon = 10^{-8}$.

Problem (p, q)	# of Iterations		CPU Time	
	MAX-MIN	CONE	MAX-MIN	CONE
(20,10)	3455	17	3.61	5.86
(40,20)	1696	15	6.90	77.25
(60,30)	1279	15	13.33	506.14
(80,40)	1183	15	25.34	2205.13
(100,50)	1073	19	40.66	8907.12
(120,60)	1017	N/A	62.90	N/A

Problem (p, q)	Memory	
	MAX-MIN	CONE
(20,10)	2.67	279
(40,20)	2.93	483
(60,30)	3.23	1338
(80,40)	3.63	4456
(100,50)	4.23	10445
(120,60)	4.98	> 16109

SUMMARY

We have shown that a (smooth) first-order method applied to a MAX-MIN reformulation of (1) substantially outperforms a first-order method applied to the corresponding dual MIN-MAX reformulation.

We have also shown that it substantially outperforms an interior-point method applied to a CP reformulation of (1).