

A Decomposition Algorithm for the Sparse Generalized Eigenvalue Problem

Ganzhao Yuan^{1,3}, Li Shen², Wei-Shi Zheng³

¹ Center for Quantum Computing, Peng Cheng Laboratory

² Tencent AI Lab

³ Sun Yat-sen University

April. 2, 2019

Outline of this talk

- The Sparse Generalized Eigenvalue Problem
- The Proposed Decomposition Algorithm
- Existing Sparse Optimization Methods
- Theoretical Analysis
- Experiments

The Sparse Generalized Eigenvalue Problem

The Sparse Generalized Eigenvalue Problem

$$\min_{\mathbf{x} \neq \mathbf{0}, \|\mathbf{x}\|_0 \leq s} f(\mathbf{x}) \triangleq \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{C} \mathbf{x}}.$$

- Statistical learning models

- ① Principle Component Analysis (PCA)

$$\min_{\mathbf{x} \neq \mathbf{0}} \frac{-\mathbf{x}^T \Sigma \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$$

- ② Fisher Discriminant Analysis (FDA)

$$\min_{\mathbf{x} \neq \mathbf{0}} \frac{-\mathbf{x}^T ((\boldsymbol{\mu}_{(1)} - \boldsymbol{\mu}_{(2)})(\boldsymbol{\mu}_{(1)} - \boldsymbol{\mu}_{(2)})^T) \mathbf{x}}{\mathbf{x}^T (\Sigma_{(1)} + \Sigma_{(2)}) \mathbf{x}}$$

- ③ Canonical Correlation Analysis (CCA)

$$\min_{\mathbf{x}} \frac{-\mathbf{x}^T \begin{pmatrix} 0 & \Sigma_{xy} \\ \Sigma_{yx} & 0 \end{pmatrix} \mathbf{x}}{\mathbf{x}^T \begin{pmatrix} \Sigma_{xx} & 0 \\ 0 & \Sigma_{yy} \end{pmatrix} \mathbf{x}}$$

- Applications: object recognition, object detection, visual tracking, pixel selection, text summarization

The Proposed Algorithm

The Concept of Block- k Optimality

- Optimization Problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x})$$

- Variable Splitting: $B \in \mathbb{N}^k$ is a subset of $\{1, 2, \dots, n\}$. $N \triangleq \{1, \dots, n\} \setminus B$, $\mathbf{x} = \mathbf{I}\mathbf{x} = (\mathbf{U}_B \mathbf{U}_B^T + \mathbf{U}_N \mathbf{U}_N^T) \mathbf{x} = \mathbf{U}_B \mathbf{x}_B + \mathbf{U}_N \mathbf{x}_N$.
- Block- k Optimal Solution

$$\mathcal{P}(B, \mathbf{x}) \triangleq \arg \min_{\mathbf{x}_B} F(\mathbf{U}_B \mathbf{x}_B + \mathbf{U}_N \mathbf{x}_N)$$

$\{\bar{\mathbf{x}}_B = \mathcal{P}(B, \bar{\mathbf{x}}) \text{ for all } |B| = k\} \Leftrightarrow \bar{\mathbf{x}}$ is the block- k optimal solution

- Block- k Optimality Measure [$\{\mathcal{B}_{(i)}\}_{i=1}^{C_n^k}$ denotes all the possible combinations]

$$\mathcal{M}(\mathbf{x}) \triangleq \frac{1}{C_n^k} \sum_{i=1}^{C_n^k} \|\mathcal{P}(\mathcal{B}_{(i)}, \mathbf{x}) - \mathbf{x}_{\mathcal{B}_{(i)}}\|_2^2$$

$$\mathcal{M}(\bar{\mathbf{x}}) = 0 \Leftrightarrow \bar{\mathbf{x}}$$
 is the block- k optimal solution

Sparse Generalized Eigenvalue Problem

Optimization Problem:

$$\min_{\mathbf{x} \neq \mathbf{0}, \|\mathbf{x}\|_0 \leq s} f(\mathbf{x}) \triangleq \frac{\frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x}}{\frac{1}{2} \mathbf{x}^T \mathbf{C} \mathbf{x}}$$

We define $h(\mathbf{x}) \triangleq \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x}$, $g(\mathbf{x}) \triangleq \frac{1}{2} \mathbf{x}^T \mathbf{C} \mathbf{x}$.

$$h(\mathbf{x}_B, \mathbf{x}_N) = \frac{1}{2} \mathbf{x}_B^T \mathbf{A}_{BB} \mathbf{x}_B + \frac{1}{2} \mathbf{x}_N^T \mathbf{A}_{NN} \mathbf{x}_N + \langle \mathbf{x}_B, \mathbf{A}_{BN} \mathbf{x}_N \rangle,$$

$$g(\mathbf{x}_B, \mathbf{x}_N) = \frac{1}{2} \mathbf{x}_B^T \mathbf{C}_{BB} \mathbf{x}_B + \frac{1}{2} \mathbf{x}_N^T \mathbf{C}_{NN} \mathbf{x}_N + \langle \mathbf{x}_B, \mathbf{C}_{BN} \mathbf{x}_N \rangle.$$

The Proposed Decomposition Algorithm

Input: k , θ , a feasible solution \mathbf{x}^0 , $t = 0$

while *not converge* **do**

S1 Find a working set B of size k . Denote $N \triangleq \{1, \dots, n\} \setminus B$.

S2 Solve the following subproblem *globally*:

$$\begin{aligned} \mathbf{x}_B^{t+1} \leftarrow \arg \min_{\mathbf{x}_B} & \frac{h(\mathbf{x}_B, \mathbf{x}_N^t) + \frac{\theta}{2} \|\mathbf{x}_B - \mathbf{x}_B^t\|_2^2}{g(\mathbf{x}_B, \mathbf{x}_N^t)} \\ \text{s.t. } & \|\mathbf{x}_B\|_0 + \|\mathbf{x}_N^t\|_0 \leq s \end{aligned} \quad (1)$$

S2 Increment t by 1

end

Algorithm 1: The Proposed Decomposition Algorithm

Remarks on the Decomposition Algorithm

- 1 A new proximal strategy (only applied to the numerator) \Rightarrow sufficient descent and global convergence
- 2 When $k = n$, the subproblem reduces to the original problem.
- 3 Finding the working set
 - 1 **Random strategy.** Select one combination (which contains k coordinates) from the whole working set $\{\mathcal{B}_{(i)}\}_{i=1}^{C_n^k}$ uniformly.
 - 2 **Swapping strategy.** Pick the top pairs of coordinates that lead to the greatest descent by measuring $\mathbf{D} \in \mathbb{R}^{|\mathcal{S}(\mathbf{x})| \times |\mathcal{Z}(\mathbf{x})|}$:

$$\mathbf{D}_{i,j} = \min_{\beta} f(\mathbf{x}^t + \beta \mathbf{e}_i - \mathbf{x}_j^t \mathbf{e}_j) - f(\mathbf{x}^t).$$

with $\mathcal{S}(\mathbf{x}) \triangleq \{i \mid \mathbf{x}_i \neq 0\}$ and $\mathcal{Z}(\mathbf{x}) \triangleq \{j \mid \mathbf{x}_j = 0\}$.

Remarks on the Decomposition Algorithm

- ④ The subproblem reduces to the following problem:

$$\min_{\mathbf{z} \in \mathbb{R}^k, \|\mathbf{z}\|_0 \leq q} p(\mathbf{z}) \triangleq \frac{\frac{1}{2} \mathbf{z}^T \bar{\mathbf{Q}} \mathbf{z} + \bar{\mathbf{p}}^T \mathbf{z} + \bar{w}}{\frac{1}{2} \mathbf{z}^T \bar{\mathbf{R}} \mathbf{z} + \bar{\mathbf{c}}^T \mathbf{z} + \bar{v}}$$

Our solution:

- We consider the following problem:

$$\min_{\mathbf{z} \in \mathbb{R}^k} p(\mathbf{z}), \text{ s.t. } \mathbf{z}_K = \mathbf{0}$$

where K has $\sum_{i=0}^q C_k^i$ possible choices.

- It reduces to the following problem:

$$\min_{\mathbf{y}} \mathcal{L}(\mathbf{y}) \triangleq \frac{\frac{1}{2} \mathbf{y}^T \mathbf{Q} \mathbf{y} + \mathbf{p}^T \mathbf{y} + w}{\frac{1}{2} \mathbf{y}^T \mathbf{R} \mathbf{y} + \mathbf{c}^T \mathbf{y} + v}$$

Solving the Quadratic Fractional Problem Globally

$$\min_{\mathbf{y}} \mathcal{L}(\mathbf{y}) \triangleq \frac{\frac{1}{2}\mathbf{y}^T \mathbf{Q} \mathbf{y} + \mathbf{p}^T \mathbf{y} + w}{\frac{1}{2}\mathbf{y}^T \mathbf{R} \mathbf{y} + \mathbf{c}^T \mathbf{y} + v}$$

$$\mathcal{J}(\alpha) = 0, \text{ with } \mathcal{J}(\alpha) \triangleq \min_{\mathbf{y}} u(\mathbf{y}) - \alpha q(\mathbf{y})$$

We have the following results.

- 1 It holds that: $\lambda_{\min}(\mathbf{Z}) \leq \min_{\mathbf{y}} \mathcal{L}(\mathbf{y}) < \lambda_{\min}(\mathbf{O})$.
- 2 The dual equivalent function $\mathcal{J}(\alpha)$ is monotonically decreasing on the range $\lambda_{\min}(\mathbf{Z}) \leq \alpha < \lambda_{\min}(\mathbf{O})$.
- 3 The optimal solution $\mathbf{y}^* = \mathbf{R}^{-1/2}(\mathbf{u}^* - \mathbf{R}^{-1/2}\mathbf{c})$, with $\mathbf{u}^* = -(\mathbf{O} - \alpha^*\mathbf{I})^{-1}\mathbf{g}$ and α^* being the unique root of the equation $\mathcal{J}(\alpha) = 0$ on the range $\lambda_{\min}(\mathbf{Z}) \leq \alpha < \lambda_{\min}(\mathbf{O})$.

Contributions of this paper:

- 1 A decomposition algorithm, it finds stronger stationary points.
- 2 Two strategies to find the working set
- 3 Two methods to solve the subproblem
- 4 A convergence analysis for 'DEC'
- 5 'DEC' consistently outperforms existing methods

Existing Sparse Optimization Methods

Optimization Problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \text{ s.t. } \|\mathbf{x}\|_0 \leq s$$

- $f(\mathbf{x})$: smooth, convex, its gradient is L -Lipschitz continuous
- Existing methods
 - 1 Relaxed approximation method
 - 2 Greedy pursuit method
 - 3 Combinatorial search method
 - 4 Gradient projection method

1 Relaxed approximation method

- convex: $f(\mathbf{x}) + \lambda \|\mathbf{x}\|_1$, $f(\mathbf{x}) + \lambda \|\mathbf{x}\|_{\text{tok}-k}$
- nonconvex: $f(\mathbf{x}) + \lambda \|\mathbf{x}\|_p$, reweighted ℓ_1 norm

⇒ our method directly controls the sparsity of the solution

2 Greedy pursuit method

- the solution MUST be initialized to zero
- $S = \emptyset$, $S = S \cup i_1$, $\min_{\mathbf{x}_S} f(\mathbf{x})$, $S = S \cup i_2$, $\min_{\mathbf{x}_S} f(\mathbf{x})$, ...

⇒ our method is a greedy coordinate descent algorithm

- ③ Combinatorial search method: global optimization methods
 - cutting plane methods
 - branch-and-cut methods

⇒ our method leverages the effectiveness of combinatorial search method methods

- ④ Gradient projection method

- $\mathbf{x}^{k+1} = \text{Proj}_s(\mathbf{x}^k - \gamma \nabla f(\mathbf{x}^k))$
- with $\text{Proj}_s(\mathbf{a}) = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{a}\|_2^2, \text{ s.t. } \|\mathbf{x}\|_0 \leq s$

⇒ our method significantly outperforms gradient projection method

Theoretical Analysis

Basic Stationary Point

$\check{\mathbf{x}}$ is called a basic stationary point if the following holds.

$\check{\mathbf{x}} = \arg \min_{\mathbf{y}} f(\mathbf{y}), \text{ s.t. } \mathbf{y}_Z = \mathbf{0}, |S| \leq s, \text{ where } Z \triangleq \{i | \check{\mathbf{x}}_i = 0\} \text{ and } S \triangleq \{i | \check{\mathbf{x}}_i \neq 0\}.$

L -Stationary Point

A solution $\check{\mathbf{x}}$ is an L -stationary point if it holds that:

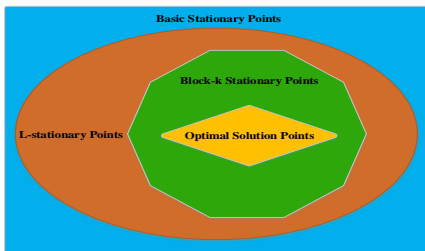
$\check{\mathbf{x}} = \text{Proj}_S(\check{\mathbf{x}} - \nabla f(\check{\mathbf{x}})/L).$

Block- k Stationary Point

A solution $\bar{\mathbf{x}}$ is a block- k stationary point if it holds that:

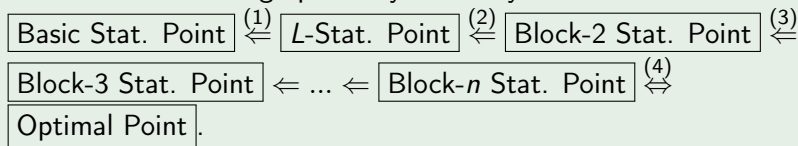
$\bar{\mathbf{x}} \in \arg \min_{\mathbf{z} \in \mathbb{R}^n} \mathcal{P}(\mathbf{z}; \bar{\mathbf{x}}, B) \triangleq \{F(\mathbf{z}), \text{ s.t. } \mathbf{z}_N = \bar{\mathbf{x}}_N\}, \forall |B| = k, N \triangleq \{1, \dots, n\} \setminus B.$

Optimality Hierarchy



Relations between the three types of stationary point.

We have the following optimality hierarchy ^a:



^aGanzhao Yuan, Li Shen, Wei-Shi Zheng. A Hybrid Method of Combinatorial Search and Coordinate Descent for Discrete Optimization. arXiv preprint, 2017. url: <https://arxiv.org/abs/1706.06493>.

A Running Example

Optimization Problems:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{x}^T \mathbf{p}, \text{ s.t. } \|\mathbf{x}\|_0 \leq s$$

$$n = 6, \mathbf{Q} = \mathbf{c} \mathbf{c}^T + \mathbf{I}, \mathbf{p} = \mathbf{1}, \mathbf{c} = [1 \ 2 \ 3 \ 4 \ 5 \ 6]^T, s = 4.$$

Number of points satisfying optimality conditions.

Basic-Stat.	L-Stat.	Block-1 Stat.	Block-2 Stat.	Block-3 Stat.	Block-4 Stat.	Block-5 Stat.	Block-6 Stat.
57	56		3	1	1	1	1

Global Convergence Properties.

Assume that the subproblem is solved globally. We have the following results.

- 1 When the random strategy is used to find the working set, we have $\lim_{t \rightarrow \infty} \mathbb{E}[\|\mathbf{x}^{t+1} - \mathbf{x}^t\|] = 0$ and Algorithm 1 converges to the block- k stationary point in expectation.
- 2 When the swapping strategy is used to find the working set with $k \geq 2$, we have $\lim_{t \rightarrow \infty} \|\mathbf{x}^{t+1} - \mathbf{x}^t\| = 0$ and Algorithm 1 converges to the block-2 stationary point deterministically.

Experiments

Sparse Generalized Eigenvalue Problem

Optimization problems:

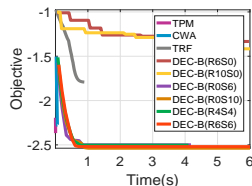
$$\min_{\mathbf{x}} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{C} \mathbf{x}}, \quad s.t. \quad \|\mathbf{x}\|_0 \leq s$$

Compared Methods:

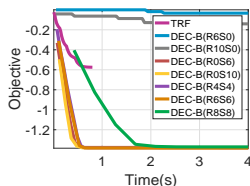
- Truncated Power Method (TPM) [Yuan & Zhang, JMLR2013]
- Coordinate-Wise Algorithm (CWA) [Beck & Vaisbourd, JOTA2016]
- Truncated Rayleigh Flow (TRF) [Tan, et al., JRSS2018]
- Quadratic Majorization Method (QMM) [Song, et al., TIP2015]
- Proposed Decomposition Method ¹ (DEC-Ri-Gj) [This Paper]

¹Selecting i coordinate using the **Random** strategy and j coordinates using the **Swapping** strategy

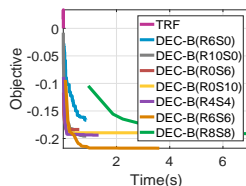
Convergence Behavior: i



(a) sparse PCA, 'rnd-2000'



(b) sparse FDA, 'rnd-2000'

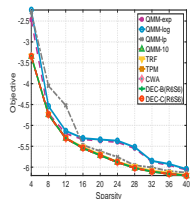


(c) sparse CCA, 'rnd-2000'

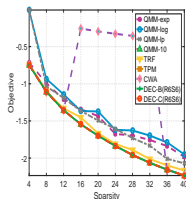
Conclusions:

- 1 {TPM,CWA,TRF} converge faster, but result in poor accuracy.
- 2 DEC-B(R10S0) achieves a lower objective value than DEC-B(R6S0). Larger k implies stronger stationary points.
- 3 The swapping strategy plays an indispensable role.

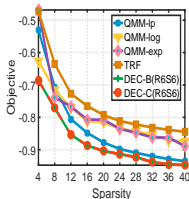
Experimental Results



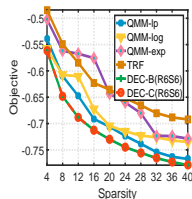
(a) a1a



(b) w1a



(c) a1a

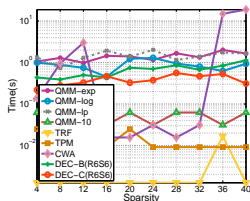


(d) w1a

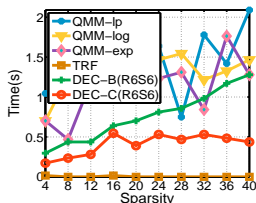
Conclusions:

- 1 CWA is not stable (much worse results on 'w1a').
- 2 DEC achieves lowest objective values.
- 3 Both DEC-B and DEC-C perform similarly.

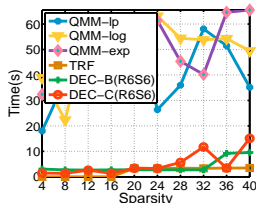
Computational Efficiency



(e) PCA



(f) FDA



(g) CCA

Conclusions:

- 1 DEC takes less than 15 seconds to converge in all our instances.
- 2 DEC is practical and it is much more efficient than QMM.

Thank You!