

Online Sparse Representation Clustering for Evolving Data Streams

—Supplementary Document

Jie Chen, *Member, IEEE*, Shengxiang Yang, *Senior Member, IEEE*, Conor Fahy, Zhu Wang, Yi-nan Guo, and Yingke Chen

I. PROOF OF THEOREM 1

In this section, we prove Theorem 1 in the paper regarding the optimization program

$$\begin{aligned} \min_{\mathbf{Z}, \mathbf{J}, \mathbf{W}} \|\mathbf{Z}\|_0 + \frac{\lambda}{2} \|\mathbf{W}^T \mathbf{X} - \mathbf{W}^T \mathbf{D} \mathbf{J}\|_F^2 \\ \text{s.t. } \mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} = \mathbf{I}_m, \mathbf{J} = \mathbf{Z} - \text{diag}(\mathbf{Z}). \end{aligned} \quad (1)$$

Given the fixed \mathbf{J}_{k+1} and \mathbf{W}_{k+1} , \mathbf{Z}_{k+1} is updated by the following objective function:

$$\mathbf{Z}_{k+1} = \min_{\mathbf{Z}_{k+1}} \frac{1}{\mu_k} \|\mathbf{Z}_{k+1}\|_0 + \frac{1}{2} \left\| \mathbf{Z}_{k+1} - \left(\mathbf{J}_{k+1} + \frac{\mathbf{Y}_k}{\mu_k} \right) \right\|_F^2. \quad (2)$$

Given a positive number $\lambda > 0$, the hard thresholding operator $\mathcal{T}_{\sqrt{\lambda}}(\mathbf{Y})$ is defined as follows [1]:

$$\mathcal{T}_{\sqrt{\lambda}}(x) = \begin{cases} 0, & \text{if } |x| \leq \sqrt{\lambda} \\ x, & \text{if } |x| > \sqrt{\lambda} \end{cases} \quad (3)$$

where $\mathbf{Y} \in \mathbb{R}^{m \times n}$ is a matrix and x represents an element of \mathbf{Y} . The closed-form solution of (2) is obtained by using the operator \mathcal{T} :

$$\mathbf{Z}_{k+1} = \mathcal{T}_{\sqrt{\frac{1}{\mu_k}}} \left(\mathbf{J}_{k+1} + \frac{\mathbf{Y}_k}{\mu_k} \right). \quad (4)$$

Theorem 1 *The convergence condition $\|\mathbf{Z}_k - \mathbf{J}_k\|_{\max} < \varepsilon$ will eventually be satisfied as k increases if ρ and μ satisfy the following conditions:*

$$\rho > 2 \quad \text{and} \quad \mu > 0$$

where k represents the number of iterations and ε is a small positive number, e.g., $\varepsilon = 10^{-4}$.

J. Chen is with the College of Computer Science, Sichuan University, Chengdu 610065, China (E-mail: chenjie2010@scu.edu.cn).

S. Yang is with the School of Computer Science and Informatics, De Montfort University, Leicester LE1 9BH, U.K. (e-mail: syang@dmu.ac.uk).

C. Fahy are with the School of Computer Science and Informatics, De Montfort University, Leicester LE1 9BH, U.K. (e-mail: conor.fahy@dmu.ac.uk).

Z. Wang is with the Law School, Sichuan University, Chengdu 610065, China (E-mail: wangzhu@scu.edu.cn).

Y. Guo is with the School of Mechanical Electronic and Information Engineering, China University of Mining and Technology (Beijing), Beijing, 100083, China. E-mail: nanfly@126.com.

Y. Chen is with the Department of Computer and Information Sciences, Northumbria University, Newcastle upon Tyne, NE1 8ST, U. K. (E-mail: yke.chen@gmail.com).

Proof According to (3), \mathbf{Z}_{k+1} has a closed-form solution in (2). Thus, we have:

$$\|\mathbf{Z}_{k+1} - \mathbf{J}_{k+1}\|_{\max} = \left\| T_{\sqrt{\frac{1}{\mu_k}}} \left(\mathbf{J}_{k+1} + \frac{\mathbf{Y}_k}{\mu_k} \right) - \mathbf{J}_{k+1} \right\|_{\max}.$$

Suppose that $\rho > 2$ and $\mu > 0$, and we obtain $\mu_k \rightarrow \infty$ when $k \rightarrow \infty$ according to $\mu_k = \rho \mu_{k-1}$. This indicates that we will obtain

$$T_{\sqrt{\frac{1}{\mu_k}}} \left(\mathbf{J}_{k+1} + \frac{\mathbf{Y}_k}{\mu_k} \right) = \mathbf{J}_{k+1} + \frac{\mathbf{Y}_k}{\mu_k}$$

as k steadily increases. Thus, we have

$$\begin{aligned} \|\mathbf{Z}_k - \mathbf{J}_k\|_{\max} - \|\mathbf{Z}_{k+1} - \mathbf{J}_{k+1}\|_{\max} \\ = \left\| \frac{\mathbf{Y}_{k-1}}{\mu_{k-1}} \right\|_{\max} - \left\| \frac{\mathbf{Y}_k}{\mu_k} \right\|_{\max} \\ = \frac{\rho \|\mathbf{Y}_{k-1}\|_{\max} - \|\mathbf{Y}_k\|_{\max}}{\mu_k} \\ = \frac{\rho \|\mathbf{Y}_{k-1}\|_{\max} - \|\mathbf{Y}_{k-1} + \mu_{k-1}(\mathbf{Z}_k - \mathbf{J}_k)\|_{\max}}{\mu_k}. \end{aligned}$$

In addition, we obtain

$$\begin{aligned} \rho \|\mathbf{Y}_{k-1}\|_{\max} - (\|\mathbf{Y}_{k-1}\|_{\max} + \|\mu_{k-1}(\mathbf{Z}_k - \mathbf{J}_k)\|_{\max}) \\ = (\rho - 1) \|\mathbf{Y}_{k-1}\|_{\max} - \mu_{k-1} \cdot \left\| \frac{\mathbf{Y}_{k-1}}{\mu_{k-1}} \right\|_{\max} \\ = (\rho - 2) \|\mathbf{Y}_{k-1}\|_{\max} \\ > 0. \end{aligned}$$

It is easy to see that the following inequality holds:

$$\begin{aligned} \|\mathbf{Y}_{k-1}\|_{\max} + \|\mu_{k-1}(\mathbf{Z}_k - \mathbf{J}_k)\|_{\max} \\ \geq \|\mathbf{Y}_{k-1} + \mu_{k-1}(\mathbf{Z}_k - \mathbf{J}_k)\|_{\max}. \end{aligned}$$

Hence,

$$\|\mathbf{Z}_k - \mathbf{J}_k\|_{\max} - \|\mathbf{Z}_{k+1} - \mathbf{J}_{k+1}\|_{\max} > 0.$$

This means there exists a certain k with two conditions, i.e., $\rho > 2$ and $\mu_1 > 0$, such that the following inequality holds:

$$\|\mathbf{Z}_{k+1} - \mathbf{J}_{k+1}\|_{\max} \leq \varepsilon$$

where $\mu = \mu_1$. Hence, convergence will eventually be achieved as k gradually increases if $\rho > 2$ and $\mu > 0$. \square

II. PROOF OF THEOREM 2

In this section, we prove Theorem 2 in the paper. We consider a general form of the $l_{2,1}$ -norm optimization problem:

$$f(\mathbf{Z}^l) = \min_{\mathbf{Z}^l} \|\mathbf{Z}^l\|_{2,1} + \frac{\beta}{2} \|\mathbf{C}^l - \mathbf{C}^l \mathbf{Z}^l\|_F^2 \quad (5)$$

s.t. $\text{diag}(\mathbf{Z}^l) = 0$

where $\beta > 0$ is a parameter. Problem (5) is a convex optimization problem. Let

$$\frac{\partial f(\mathbf{Z}^l)}{\partial \mathbf{W}} = 0 \quad (6)$$

and we have

$$\mathbf{Z}^l = \left(\frac{1}{\beta} \boldsymbol{\Sigma} + \mathbf{C}^{lT} \mathbf{C}^l \right)^{-1} \mathbf{C}^{lT} \mathbf{C}^l \quad (7)$$

where $\mathbf{Z}^l = [\mathbf{z}_1^r, \mathbf{z}_2^r, \dots, \mathbf{z}_i^r, \dots, \mathbf{z}_{n_1}^r]^T$ and $\boldsymbol{\Sigma} \in \mathbb{R}^{n_1 \times n_1}$ is a diagonal matrix whose diagonal entries are given by $\frac{1}{\|\mathbf{z}_i^r\|_2}$.

To prove Theorem 2, we need to prove Lemma 1.

Lemma 1 For two matrices $\mathbf{A} \in \mathbb{R}^{d \times m}$ and $\mathbf{B} \in \mathbb{R}^{m \times n}$, the following inequality holds:

$$\|\mathbf{AB}\|_F^2 \leq \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2.$$

Proof First, we have

$$\|\mathbf{AB}\|_F^2 = \text{tr}(\mathbf{A}^T \mathbf{A} \mathbf{B} \mathbf{B}^T)$$

by the definition of the trace function.

Second, we want to prove that

$$\text{tr}(\mathbf{A}^T \mathbf{A} \mathbf{B} \mathbf{B}^T) \leq \text{tr}(\mathbf{A}^T \mathbf{A}) \text{tr}(\mathbf{B} \mathbf{B}^T) \quad (8)$$

which implies that $\|\mathbf{AB}\|_F^2 \leq \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2$. It is easy to see that $\mathbf{A}^T \mathbf{A}$ and $\mathbf{B} \mathbf{B}^T$ are positive semidefinite and symmetric matrices. Using the singular value decomposition (SVD) results of $\mathbf{A}^T \mathbf{A}$ and $\mathbf{B} \mathbf{B}^T$, we obtain

$$\begin{aligned} \mathbf{A}^T \mathbf{A} &= \mathbf{U}_A \boldsymbol{\Sigma}_A \mathbf{U}_A^T, \\ \mathbf{B} \mathbf{B}^T &= \mathbf{U}_B \boldsymbol{\Sigma}_B \mathbf{U}_B^T, \end{aligned} \quad (9)$$

where \mathbf{U}_A and \mathbf{U}_B are unitary matrices, and $\boldsymbol{\Sigma}_A$ and $\boldsymbol{\Sigma}_B$ are diagonal matrices whose diagonal elements are singular values of $\mathbf{A}^T \mathbf{A}$ and $\mathbf{B} \mathbf{B}^T$, respectively. The singular values of $\mathbf{A}^T \mathbf{A}$ and $\mathbf{B} \mathbf{B}^T$ are all nonnegative. Then,

$$\begin{aligned} \text{tr}(\mathbf{A}^T \mathbf{A} \mathbf{B} \mathbf{B}^T) &= \text{tr}(\mathbf{U}_A \boldsymbol{\Sigma}_A \mathbf{U}_A^T \mathbf{U}_B \boldsymbol{\Sigma}_B \mathbf{U}_B^T) \\ &\leq \text{tr}(\mathbf{U}_B^T \mathbf{U}_A \boldsymbol{\Sigma}_A \mathbf{U}_A^T \mathbf{U}_B) \text{tr}(\boldsymbol{\Sigma}_B) \\ &= \text{tr}(\mathbf{U}_A \boldsymbol{\Sigma}_A \mathbf{U}_A^T) \text{tr}(\mathbf{U}_B \boldsymbol{\Sigma}_B \mathbf{U}_B^T) \\ &= \text{tr}(\mathbf{A}^T \mathbf{A}) \text{tr}(\mathbf{B} \mathbf{B}^T). \end{aligned} \quad (10)$$

Hence,

$$\|\mathbf{AB}\|_F^2 \leq \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2. \quad \square$$

The proof of Theorem 2 is motivated by a feature selection method [2].

Theorem 2 The objective value of Equation (7) will monotonically decrease until convergence to the global optimum of Problem (5).

Proof Suppose that \mathbf{Z}_{t+1}^l is the global optimal solution to Problem (5), i.e.,

$$\mathbf{Z}_{t+1}^l = \arg \min_{\mathbf{Z}^l \text{ diag}(\mathbf{Z}^l)=0} \|\mathbf{Z}^l\|_{2,1} + \frac{\beta}{2} \|\mathbf{C}^l - \mathbf{C}^l \mathbf{Z}^l\|_F^2.$$

Problem (5) is a convex optimization problem, which indicates that

$$\begin{aligned} \|\mathbf{Z}_{t+1}^l\|_{2,1} + \frac{1}{\beta} \|\mathbf{C}^l - \mathbf{C}^l \mathbf{Z}_{t+1}^l\|_F^2 \\ \leq \|\mathbf{Z}_{t+1}^l\|_{2,1} + \frac{1}{\beta} \|\mathbf{C}^l - \mathbf{C}^l \mathbf{Z}_t^l\|_F^2. \end{aligned}$$

Thus,

$$\|\mathbf{C}^l - \mathbf{C}^l \mathbf{Z}_{t+1}^l\|_F^2 \leq \|\mathbf{C}^l - \mathbf{C}^l \mathbf{Z}_t^l\|_F^2.$$

According to Lemma 1, we have

$$\|\mathbf{I} - \mathbf{Z}_{t+1}^l\|_F^2 \leq \|\mathbf{I} - \mathbf{Z}_t^l\|_F^2$$

where \mathbf{I} is an identity of size $n_l \times n_l$. Then, we have the following inequality:

$$\text{tr}(\mathbf{Z}_{t+1}^l (\mathbf{Z}_{t+1}^l)^T - \mathbf{Z}_t^l (\mathbf{Z}_t^l)^T) \leq \text{tr}(2(\mathbf{Z}_{t+1}^l - \mathbf{Z}_t^l)).$$

Using the constraint $\text{diag}(\mathbf{Z}^l) = 0$ in Problem (5), we obtain

$$\text{tr}(\mathbf{Z}_{t+1}^l (\mathbf{Z}_{t+1}^l)^T - \mathbf{Z}_t^l (\mathbf{Z}_t^l)^T) \leq 0.$$

Then

$$\sum_{i=1}^n \left\| (\mathbf{z}^i)_{t+1}^l \right\|_2^2 \leq \sum_{i=1}^n \left\| (\mathbf{z}^i)_t^l \right\|_2^2$$

where $(\mathbf{z}^i)_{t+1}^l$ and $(\mathbf{z}^i)_t^l$ are the i -th row vectors of \mathbf{Z}_{t+1}^l and \mathbf{Z}_t^l , respectively. Hence,

$$\|\mathbf{Z}_{t+1}^l\|_{2,1} \leq \|\mathbf{Z}_t^l\|_{2,1}.$$

This means that the objective value of Equation (7) will monotonically decrease at each iteration. At the $(t+1)$ -th iteration, Equation (7) holds for given \mathbf{Z}_{t+1}^l and $\boldsymbol{\Sigma}_{t+1}^l$. Consequently, the objective value of Equation (7) will converge to the global optimum of Problem (5). \square

REFERENCES

- [1] T. Blumensath and M. E. Davies, "Iterative thresholding for sparse approximations," *J. Fourier Anal. Appl.*, vol. 14, no. 5, pp. 629–654, Sept. 2008.
- [2] F. Nie, H. Huang, X. Cai, and C. Ding, "Efficient and robust feature selection via joint $l_{2,1}$ -norms minimization," in *Adv. Neural. Inf. Process. Syst.*, Vancouver, British Columbia, Canada, Dec. 2010, pp. 1813–1821.