

# Supplementary materials of “NLRRC: a novel clustering method of jointing non-negative LRR and random walk graph regularized NMF for single-cell type identification”

Juan Wang, Lin-Ping Wang, Shasha Yuan, Feng Li, Jin-Xing Liu, Member, IEEE, Jun-Liang Shang, Member, IEEE

## Support materials

### 1. The optimization of NLRRC

The objective function of NLRRC is not convex. The multiplication update rule is used to iteratively update the three variables  $\mathbf{Z}$ ,  $\mathbf{B}$ , and  $\mathbf{C}$  in the objective function.

First, the auxiliary variable  $\mathbf{S}$  is introduced to discretize the objective function of NLRRC.

Then, Lagrange multipliers  $\phi = [\phi_{ik}]$  and  $\varphi = [\varphi_{kj}]$  are introduced to constrain  $\mathbf{B} \geq \mathbf{0}$  and  $\mathbf{C} \geq \mathbf{0}$ , respectively. The objective function of NLRRC is rewritten as follows:

$$\begin{aligned} L(\mathbf{Z}, \mathbf{S}, \mathbf{B}, \mathbf{C}) = & \|\mathbf{X} - \mathbf{XS}\|_F^2 + \lambda \|\mathbf{Z}\|_* + \|\mathbf{S} - \mathbf{BC}\|_F^2 \\ & + \beta \text{tr}(\mathbf{CL}_{RWG} \mathbf{C}^T) + \mu (\|\mathbf{Z} - \mathbf{S}\|_F^2 + \|\mathbf{C}^T \mathbf{C} - \mathbf{I}\|_F^2) \\ & + \langle \mathbf{C}_1, \mathbf{Z} - \mathbf{S} \rangle + \langle \mathbf{C}_2, \mathbf{C}^T \mathbf{C} - \mathbf{I} \rangle + \text{tr}(\phi \mathbf{B}) + \text{tr}(\varphi \mathbf{C}), \end{aligned} \quad (1)$$

where  $\mu$  represents the penalty parameter,  $\mathbf{C}_1 \in \mathbf{R}^{n \times n}$  and  $\mathbf{C}_2 \in \mathbf{R}^{n \times n}$  are Lagrange multipliers, respectively.  $\langle \mathbf{A}, \mathbf{B} \rangle$  represents the inner product of  $\mathbf{A}$  and  $\mathbf{B}$ .

Finally, the alternating direction method of multiplier (ADMM) is used to iteratively process each variable to find the local optimal solution. The specific update process for each variable is as follows:

Step 1: Once  $\mathbf{S}$ ,  $\mathbf{B}$ , and  $\mathbf{C}$  are fixed, update  $\mathbf{Z}$ . The variable  $\mathbf{Z}$  is written as

$$L(\mathbf{Z}) = \lambda \|\mathbf{Z}\|_* + \mu \|\mathbf{Z} - \mathbf{S}\|_F^2 + \langle \mathbf{C}_1, \mathbf{Z} - \mathbf{S} \rangle, \quad (2)$$

Equation (2) is a quadratic optimization problem with low-rank constraints. Equation (2) is solved using the soft-threshold method to obtain the closed-form solution shown below.

$$\mathbf{Z}^{k+1} = \text{soft}_{\lambda, \frac{1}{\mu}}(\mathbf{S}^k - \frac{\mathbf{C}_1^k}{\mu}), \quad (3)$$

where  $\text{soft}_{\lambda, \frac{1}{\mu}}(\cdot)$  represents the soft-threshold operator.

Step 2: When  $\mathbf{Z}$ ,  $\mathbf{B}$ , and  $\mathbf{C}$  are fixed, update  $\mathbf{S}$ .  $\mathbf{S}$  can be expressed mathematically as

$$\begin{aligned} L(\mathbf{S}) = & \|\mathbf{X} - \mathbf{XS}\|_F^2 + \|\mathbf{S} - \mathbf{BC}\|_F^2 \\ & + \mu \|\mathbf{Z} - \mathbf{S}\|_F^2 + \langle \mathbf{C}_1, \mathbf{Z} - \mathbf{S} \rangle. \end{aligned} \quad (4)$$

Equation (4) is optimized as follows:

$$\mathbf{S}^{k+1} = (\mathbf{X}^T \mathbf{X} + \mathbf{I} + \mu \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{X} + \mathbf{B}^k \mathbf{C}^k + \mu \mathbf{Z}^{k+1} + \mathbf{C}_1^k). \quad (5)$$

then, we set  $\mathbf{S}^{k+1} = ((\mathbf{S}^{k+1})^T + \mathbf{S}^{k+1})/2$  to ensure the symmetry of similar matrices.

Step 3: Update  $\mathbf{B}$  and  $\mathbf{C}$ .  $L$  to  $\mathbf{B}$  and  $\mathbf{C}$  derivative is

$$\frac{\partial L}{\partial \mathbf{B}^k} = \mathbf{B}^k \mathbf{C}^k \mathbf{C}^{kT} - \mathbf{S}^{k+1} \mathbf{C}^{kT} + \phi, \quad (6)$$

$$\frac{\partial L}{\partial \mathbf{C}^k} = \mathbf{B}^{kT} \mathbf{B}^k \mathbf{C}^k + \mathbf{B}^k \mathbf{C}_2 + \beta \mathbf{C}^k \mathbf{L} - \mathbf{B}^{kT} \mathbf{S}^{k+1} + \varphi, \quad (7)$$

The iteration formula of objective functions (6) and (7) is expressed using the Karush-Kuhn-Tucker (KKT) [1] under the conditions of  $\phi_{ik} w_{ik} = 0$  and  $\varphi_{kj} h_{kj} = 0$ , respectively, as

$$\mathbf{B}_{ik}^{k+1} \leftarrow \mathbf{B}_{ik}^k \frac{(\mathbf{S}^{k+1} (\mathbf{C}^k)^T)_{ik}}{(\mathbf{B}^k \mathbf{C}^k (\mathbf{C}^k)^T)_{ik}}, \quad (8)$$

$$\mathbf{C}_{kj}^{k+1} \leftarrow \mathbf{C}_{kj}^k \frac{((\mathbf{B}^{k+1})^T \mathbf{S} + \beta \mathbf{C}^k \mathbf{W})_{kj}}{((\mathbf{B}^{k+1})^T \mathbf{B}^{k+1} \mathbf{C}^k + \mathbf{C}^k \mathbf{C}_2 + \beta \mathbf{C}^k \mathbf{D})_{kj}}, \quad (9)$$

where  $\mathbf{W}$  is the symmetric weight matrix and  $\mathbf{D}$  is the degree matrix of the edge.

Step 4: Update  $\mathbf{C}_1$  and  $\mathbf{C}_2$

$$\mathbf{C}_1^{k+1} = \mathbf{C}_1^k + \mu (\mathbf{Z}^{k+1} - \mathbf{S}^{k+1}), \quad (10)$$

$$\mathbf{C}_2^{k+1} = \mathbf{C}_2^k + \mu ((\mathbf{C}^{k+1})^T \mathbf{C}^{k+1} - \mathbf{I}). \quad (11)$$

Repeat the preceding steps until the set constraints or the maximum number of iterations are met. Algorithm 1 depicts the specific update process of LRRNC.

## 2. Time complexity analysis

**Algorithm 1** The algorithm for LRRNC.

---

<b>Input:</b> Gene representation matrix $\mathbf{X}$ , Parameter $k$ , $\lambda$ and $\beta$ .
<b>Initialization:</b> Initialize $\mathbf{Z}, \mathbf{S}, \mathbf{B}, \mathbf{C}$ , maximum number of iterations and stop error. $\mathbf{Z}=\mathbf{0}, \mathbf{S}=\mathbf{0}, \mathbf{B}=\mathbf{0}, \mathbf{C}=\mathbf{0}, \mu=10^{-4}$ , $\max Iter=100, tol_1=10^{-5}, tol_2=10^{-5}$ .
<b>While not convergence do</b> 1) Update the variable $\mathbf{Z}$ with Equation (3), 2) Update the variable $\mathbf{S}$ with Equation (5), 3) Update the variable $\mathbf{B}$ and $\mathbf{C}$ with Equation (8) and Equation (9) respectively. 4) Update the Lagrange multiplier $\mathbf{C}_1$ and $\mathbf{C}_2$ with Equation (10) and Equation (11), respectively. <b>Until the maximum iteration is reached or satisfies the following convergence conditions:</b> $\max(\ \mathbf{X}-\mathbf{XZ}\ _F^2, \ \mathbf{Z}-\mathbf{BC}\ _F^2) / \max(\ \mathbf{XZ}\ _F^2, \ \mathbf{BC}\ _F^2) < tol_1$ , $\max(\ \mathbf{Z}^{k+1}-\mathbf{Z}^k\ _F^2, \ \mathbf{S}^{k+1}-\mathbf{S}^k\ _F^2, \ \mathbf{B}^{k+1}-\mathbf{B}^k\ _F^2, \ \mathbf{C}^{k+1}-\mathbf{C}^k\ _F^2) / \max(\ \mathbf{X}\ _F^2) < tol_2$ . <b>End while</b>
<b>Output:</b> Optimal matrix $\mathbf{C}$

---

The computational complexity of the NLRRC method is primarily caused by the updating of variables  $\mathbf{Z}, \mathbf{S}, \mathbf{B}$  and  $\mathbf{C}$ . The symbol  $O$  represents the computational complexity. Let  $n$  represent the number of cell samples. The  $m$  represents the feature numbers of each sample. The  $k$  is the number of clusters. The  $t$  represents the number of model iterations. **Algorithm 1** depicts the four main steps of NLRRC, with the first, second, and third steps requiring the highest calculation cost. The calculation cost of the first step is primarily generated by SVD solving the soft threshold function of  $\mathbf{Z}$ , and the calculation complexity is  $O(n^3)$ . The second step is  $O(n^3)$ , which is primarily generated by solving the similarity matrix. The time complexity generated by NMF is  $O(n^2 k)$  in the third step. Additionally,  $O(mn^2)$  is required for data graph construction. The computational complexity for RWGR is approximately  $O(n^2)$ . The total computational complexity of NLRRC is  $O(2n^3 t + n^2 k t + mn^2 + n^2)$ , if the model is iterated  $t$  times.

### 3. Evaluation metrics

Normalized Mutual Information (NMI) [2] and the Adjusted Rand Index (ARI) [3] are used to evaluate clustering performance. The NMI theoretically explains the consistency between the predicted cluster label  $\mathbf{Y}=\{Y_1, Y_2, \dots, Y_k\}$  and the actual cluster label  $\mathbf{T}=\{T_1, T_2, \dots, T_k\}$ . The following is the definition of NMI.

$$NMI(\mathbf{Y}, \mathbf{T}) = \frac{MI(\mathbf{Y}, \mathbf{T})}{[E(\mathbf{Y}) + E(\mathbf{T})]/2}, \quad (12)$$

where  $MI(\cdot, \cdot)$  and  $E(\cdot)$  represent mutual information and information entropy, respectively.

ARI evaluates performance by calculating similarity between the predicted label  $\mathbf{Y}$  and the actual label  $\mathbf{T}$ . ARI is

mathematically expressed as follows:

$$ARI(\mathbf{Y}, \mathbf{T}) = \frac{\binom{n}{2} (a_{\mathbf{YT}} + a) - \left[ \frac{(a_{\mathbf{YT}} + a_{\mathbf{T}})(a_{\mathbf{YT}} + a_{\mathbf{Y}})}{+ (a_{\mathbf{Y}} + a)(a_{\mathbf{T}} + a)} \right]}{\binom{n}{2} - \left[ \frac{(a_{\mathbf{YT}} + a_{\mathbf{T}})(a_{\mathbf{YT}} + a_{\mathbf{Y}})}{+ (a_{\mathbf{Y}} + a)(a_{\mathbf{T}} + a)} \right]}, \quad (13)$$

where  $a_{\mathbf{YT}}$  denotes the number of cells shared by prediction label  $\mathbf{Y}$  and real label  $\mathbf{T}$ .  $a_{\mathbf{T}}$  and  $a_{\mathbf{Y}}$  are the numbers of single-cell in real label  $\mathbf{T}$  and prediction label  $\mathbf{Y}$ , respectively.  $a$  denotes the number of cells that are not shared by real label  $\mathbf{T}$  and prediction label  $\mathbf{Y}$ .

The value range of NMI is [0,1], and the value range of ARI is [-1,1]. The higher the values of NMI and ARI are, the better the consistency between the predicted and actual labels.

### 4. Convergence

The convergence of the NLRRC is described in detail in this section. The results are displayed in Fig. 1. The vertical axis  $Y$  represents the iteration error, and the horizontal axis  $X$  represent the number of iterations. After the parameters of the model are fixed, count the error of each iteration until the maximum number of iterations or the convergence conditions set by us are met. In Fig. 1, with the increase of iteration times, we found that NLRRC converges to a lower error and tends to be stable. As a result, the NLRRC method converges well.

### 5. Selecting the number of clusters

Estimating the number of clusters is one of the most important clustering steps. The Gap Statistics [4] is a reliable method for predicting the number of clusters. We use the Gap Statistic to determine the number of clusters, as shown in TABLE I. The number of correctly predicted clusters is shown in bold. The following is the operating principle of the Gap Statistic.

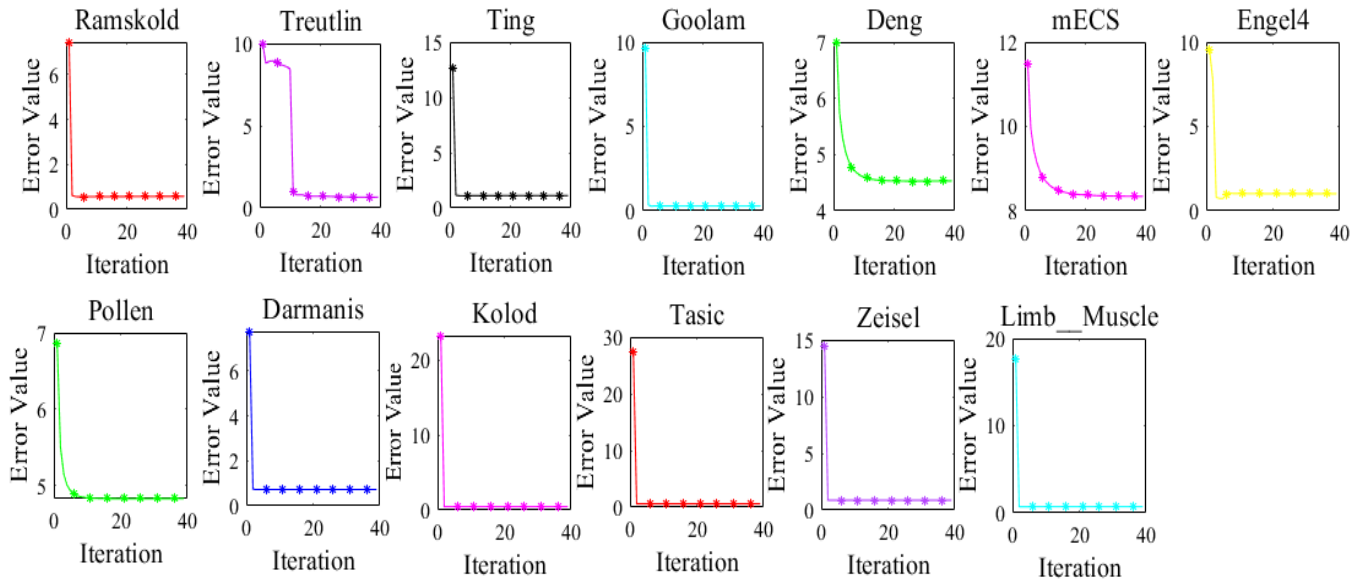


Fig. 1. Convergence curves of NLRRC on scRNA-seq datasets.

TABLE I  
EVALUATE CLUSTER NUMBER  $k$

Datasets	Ramskold	Treutlein	Ting	Goolam	Deng	mECS	Engel4	Pollen	Darmanis	Kolod	Tasic	Zeisel	Limb_Muscle
Real cluster number	7	5	5	5	7	3	4	11	8	3	48	9	6
Predict cluster number	6	5	5	5	7	3	4	12	9	3	48	9	7

Let  $Gap_{(k)}$  represent the difference between the expected value of the random sample and the Euclidean distance between the observed sample points in the clusters. If the optimal cluster number is  $k$ , the sum of Euclidean distances between observation samples within each cluster is the smallest and the value of  $Gap_{(k)}$  is the largest. The Gap Statistics employs this principle to determine the optimal cluster number by searching for the  $k$  that maximizes  $Gap_{(k)}$ .

In TABLE I, we discover that the Gap Statistics accurately predict the number of clusters in nine datasets, including Treutlein, Ting, Goolam, Deng, mECS, Engel4, Kolod, Tasic, and Zeisel.

- [4] R. Tibshirani, and W. Hastie, "Estimating the number of clusters in a data set via the gap statistic," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 63, no. 2, pp. 411-423, 2001.

## REFERENCES

- [1] A. Dreves, et al., "On the solution of the KKT conditions of generalized Nash equilibrium problems," *Siam J. Optim.*, vol. 21, no. 3, pp. 1082-1108, 2011.
- [2] A. Strehl, and J. Ghosh, "Cluster Ensembles - A Knowledge Reuse Framework for Combining Multiple Partitions," *J Mach Learn Res.*, vol. 3, no. 3, pp. 583-617, 2002.
- [3] S. Wagner, and D. Wagner, "Comparing Clusterings - An Overview," *Karlsruhe: Universität Karlsruhe, Fakultät für Informatik*, no.1, pp.1-19, 2007.