

# CELL SUBCLASS IDENTIFICATION IN SINGLE-CELL RNA-SEQUENCING DATA USING ORTHOGONAL NONNEGATIVE MATRIX FACTORIZATION

Shuai Wang<sup>†</sup>, Peng Wu<sup>‡†</sup>, Manqi Zhou<sup>†</sup>, Tsung-Hui Chang<sup>†\*</sup> and Song Wu<sup>‡</sup>

<sup>†</sup> The Chinese Univ. of Hong Kong, Shenzhen, <sup>\*</sup>Shenzhen Research Institute of Big Data, <sup>‡</sup>Institute of Urological Surgery of Shenzhen Univ.

## I. Introduction

Identification of cell subclasses using scRNA-Seq data is of paramount importance since it uncovers the hidden biological processes within the cell population. Our contributions are as follows:

- Propose the use of orthogonally constrained NMF (ONMF) model and a computationally efficient algorithm based on variable splitting and ADMM for subclass identification.
- Obtain promising results in identifying cell subclasses and detecting key (biologically meaningful) genes on real-world data.

## II. Problem Formulation

Consider a scRNA-Seq dataset consisting of the expression levels of  $M$  genes of  $N$  cell samples, denoted by matrix  $X \in \mathbb{R}^{M \times N}$ .

### The NMF Model

$$\min_{W, H} \|X - WH\|_F^2 \quad \text{s.t. } W \geq 0, H \geq 0, \quad (1)$$

where  $\|\cdot\|_F$  denotes the Frobenius norm and  $W \geq 0$  ( $H \geq 0$ ) means that all elements of  $W$  ( $H$ ) are nonnegative.

- NMF has been shown to be powerful in detecting subclasses among cell samples.

- NMF may still fail in clustering some datasets with heterogeneous structures.

### The ONMF Model

$$\min_{W, H} \|X - WH\|_F^2 \quad (2)$$

$$\text{s.t. } W \geq 0, H \geq 0, HH^T = I_K. \quad (3)$$

- Orthogonally constrained NMF (ONMF) is closely related to K-means clustering.
- The orthogonality and non-negativity constraint enforces  $H$  to be a cluster indicator matrix.

$$[H]_{k,n} = \begin{cases} c \neq 0 & \text{if cell } n \text{ belongs to cluster } k, \\ 0 & \text{otherwise} \end{cases}$$

- The ONMF problem is more challenging to solve with orthogonality constraint.

## III. Proposed Algorithm

### Variable Splitting

$$\min_{W, H, S, P, Y} \|X - WH\|_F^2, \quad (4)$$

$$\text{s.t. } W = S, H = P, H = Y, \quad (5)$$

$$S \geq 0, P \geq 0, YY^T = I_K. \quad (6)$$

New variables ( $S, P, Y$ ) split the non-negativity or orthogonality constraint from ( $W, H, H$ ).

### ADMM Based Updates

$$\begin{aligned} \mathcal{L}_a(W, H, S, P, Y, \Lambda) = & \frac{1}{2} \|X - WH\|_F^2 \\ & + \text{Tr}(\Lambda_1^T (W - S)) + \frac{\rho_1}{2} \|W - S\|_F^2 + \text{Tr}(\Lambda_2^T (H - P)) \\ & + \frac{\rho_2}{2} \|H - P\|_F^2 + \text{Tr}(\Lambda_3^T (H - Y)) + \frac{\rho_3}{2} \|H - Y\|_F^2, \end{aligned}$$

At each iteration  $r$ , we update as follows:

$$W^{r+1} \leftarrow \arg \min_W \mathcal{L}_a(W, H^r, S^r, P^r, Y^r, \Lambda^r),$$

$$H^{r+1} \leftarrow \arg \min_H \mathcal{L}_a(W^{r+1}, H, S^r, P^r, Y^r, \Lambda^r),$$

$$S^{r+1} \leftarrow \arg \min_{S \geq 0} \mathcal{L}_a(W^{r+1}, H^{r+1}, S, P^r, Y^r, \Lambda^r),$$

$$P^{r+1} \leftarrow \arg \min_{P \geq 0} \mathcal{L}_a(W^{r+1}, H^{r+1}, S^{r+1}, P, Y^r, \Lambda^r),$$

$$Y^{r+1} \leftarrow \arg \min_{YY^T = I_K} \mathcal{L}_a(W^{r+1}, H^{r+1}, S^{r+1}, P^{r+1}, Y, \Lambda^r),$$

$$\Lambda_1^{r+1} \leftarrow \Lambda_1^r + \rho_1 (W^{r+1} - S^{r+1}),$$

$$\Lambda_2^{r+1} \leftarrow \Lambda_2^r + \rho_2 (H^{r+1} - P^{r+1}),$$

$$\Lambda_3^{r+1} \leftarrow \Lambda_3^r + \rho_3 (H^{r+1} - Y^{r+1}).$$

Each of above steps has closed-form solution.

$$W^{r+1} \leftarrow (X(H^r)^T + \rho_1 S^r - \Lambda_1^r) [H^r (H^r)^T + \rho_1 I_K]^{-1}$$

$$\begin{aligned} H^{r+1} \leftarrow & (X^T W^{r+1} + \rho_2 P^r - \Lambda_2^r + \rho_3 Y^r - \Lambda_3^r) \\ & \times [(W^{r+1})^T W^{r+1} + (\rho_2 + \rho_3) I_K]^{-1}. \end{aligned}$$

$$S^{r+1} \leftarrow \max(W^{r+1} + \Lambda_1^r / \rho_1, 0),$$

$$P^{r+1} \leftarrow \max(H^{r+1} + \Lambda_2^r / \rho_2, 0),$$

$$Y^{r+1} \leftarrow VU^T \quad (SVD \text{ of } (H^{r+1} + \frac{1}{\rho_3} \Lambda_3^r) = V\Sigma U^T)$$

### Subclass Identification & Key Gene Extraction

- Apply k-means on  $H$  with an initial subclass association from  $H$  to get clustering result.

- Adopt the scoring scheme [1] to obtain the significance of each gene from rows of  $W$  and key genes are those top ranked.

## IV. Numerical Result I

### Target Datasets

Dataset	Samples	Genes	Clusters
1 Mouse Embryonic Fibroblasts	405	12117	5
2 Bladder Cancer	121	23048	4

Figure 1: ScRNA-Seq datasets

### Algorithm Convergence Performance

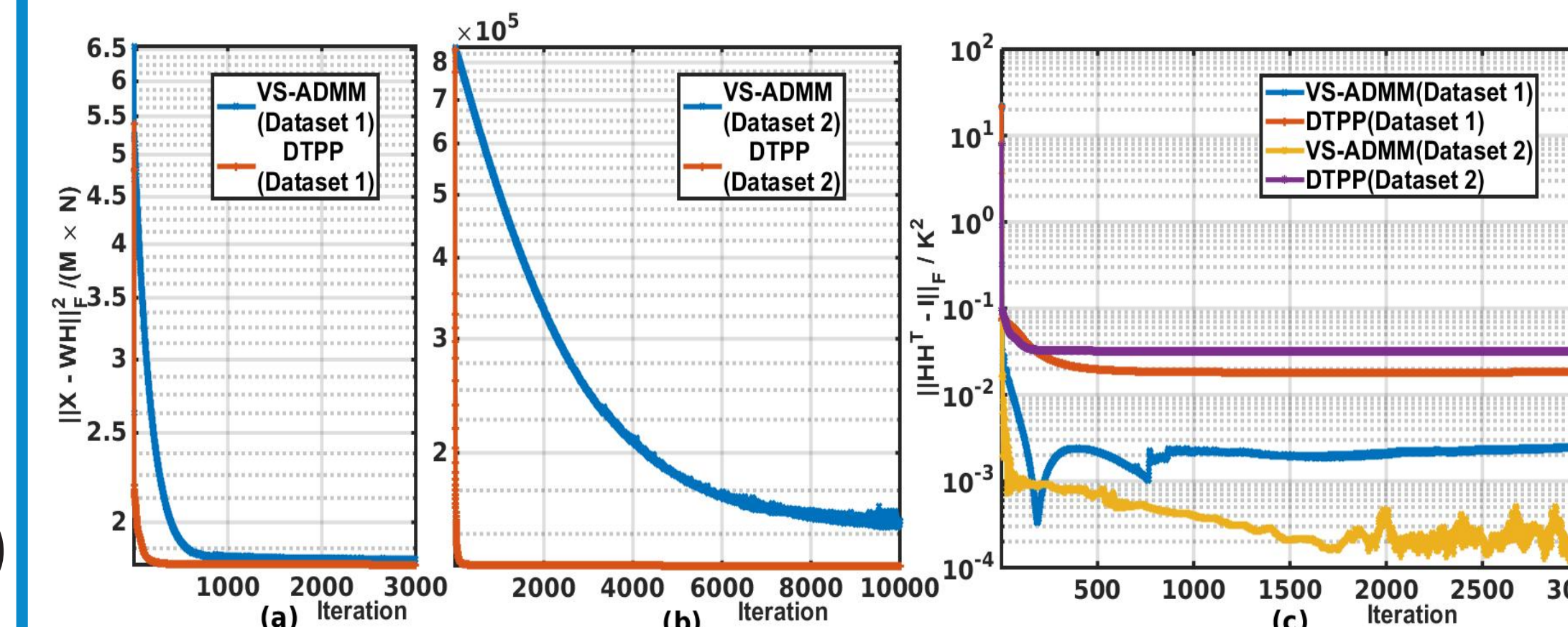


Figure 2: The convergence of objective value and feasibility of orthogonality constraint of DTPP and proposed VS-ADMM algorithms applied to Datasets 1 and 2.

### Biological Significance Analysis

Cluster	Biological Pathway	Genes	FDR q-value
1 (57 genes)	LIPID_TRANSPORTER_ACTIVITY	3	1.89E-1
	ENDOMEMBRANE_SYSTEM_ORGANIZATION	4	2.07E-1
	CELLULAR_LIPID_METABOLIC_PROCESS	5	2.07E-1
	NEGATIVE_REGULATION_OF_LIPASE_ACTIVITY	2	3.03E-2
	REGULATION_OF_LIPASE_ACTIVITY	2	4.87E-1
2 (24 genes)	HISTONE_DEMETHYLASE_ACTIVITY	2	1.27E-1
	HISTONE_BINDING	3	1.27E-1
	IMMUNE_SYSTEM	33	1.85E-3
3 (58 genes)	CELL_CYCLE	19	3.46E-3
	MISMATCH_REPAIR	4	1.52E-2
	PPAR_SIGNALING_PATHWAY	6	1.90E-2
	CELL_CYCLE_CHECKPOINTS	8	1.90E-2
	REGULATION_OF_ACTIN_CYTOSKELETON	10	3.99E-2
	CELL_CYCLE_CHECKPOINTS	8	1.90E-2
	CELL_CYCLE_CHECKPOINTS	8	1.90E-2
	CELL_CYCLE_CHECKPOINTS	8	1.90E-2
4 (861 genes)	CELL_CYCLE_CHECKPOINTS	8	1.90E-2
	CELL_CYCLE_CHECKPOINTS	8	1.90E-2
	CELL_CYCLE_CHECKPOINTS	8	1.90E-2
	CELL_CYCLE_CHECKPOINTS	8	1.90E-2

## V. Reference

- H. Kim et al., "Sparse non-negative matrix factorization via alternating non-negative-constrained least squares for microarray data analysis," *Bioinformatics*, vol. 23, no. 12, pp. 1495-1502, May 2007.
- D. D. Lee et al., "Algorithms for non-negative matrix factorization," in *Proc. NIPS*, Denver, CO, USA, Dec. 2000, pp. 556-562.
- C. Ding et al., "Orthogonal nonnegative matrix t-factorizations for clustering," in *Proc. KDD*, Philadelphia, PA, USA, Aug. 2006, pp. 20-23.

## IV. Numerical Result II

### Subclass Identification Performance

Table 1: Clustering Performance of Different Methods for Dataset 1

	Purity	Rand Index	Silhouette
K-means	0.708	0.427	0.060
NMF (Euclidean) [2]	0.731	0.483	0.538
NMF (KL) [2]	0.742	0.489	0.616
DTPP [3]	0.741	0.491	0.680
<b>Proposed VS-ADMM</b>	<b>0.749</b>	<b>0.506</b>	<b>0.803</b>

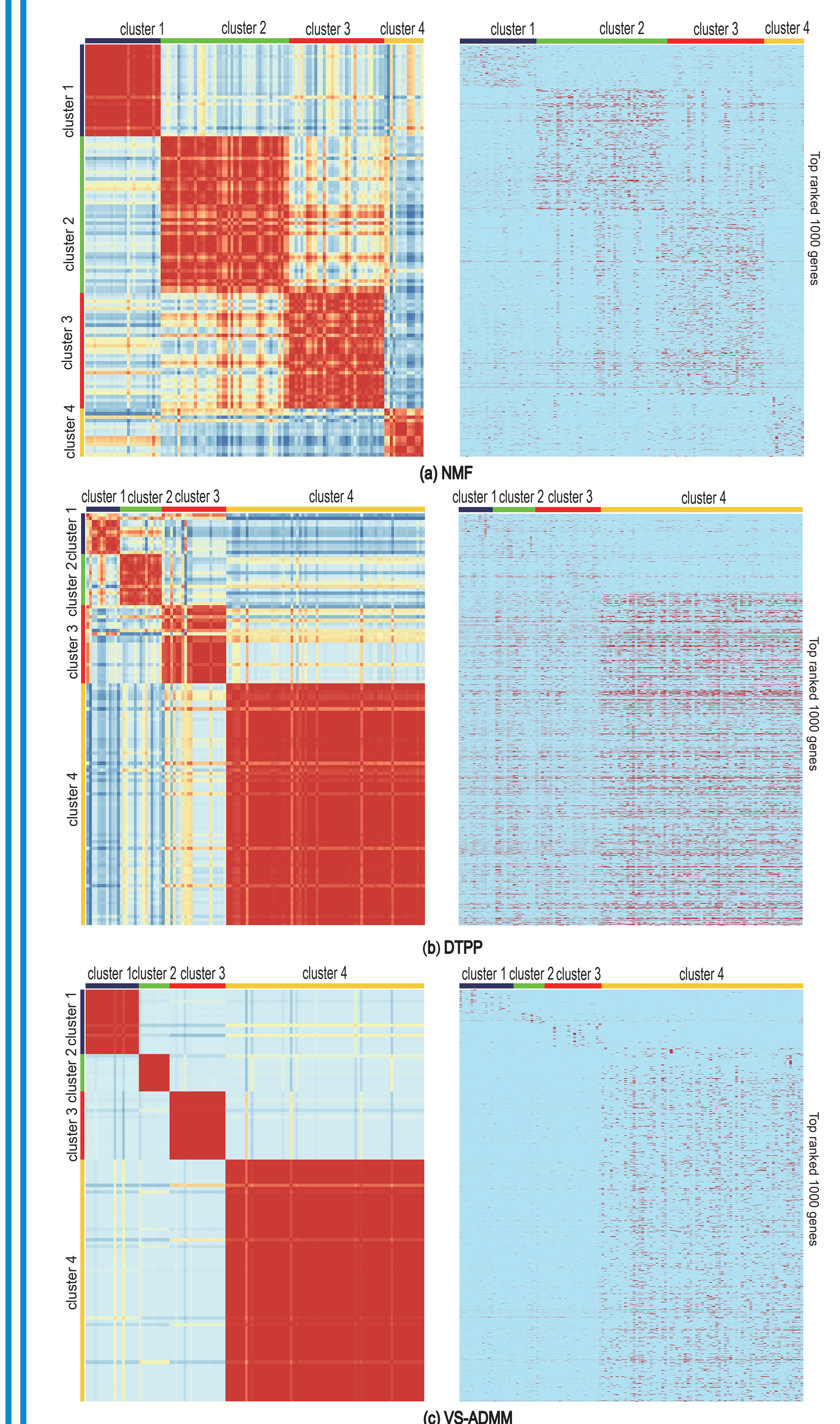


Figure 3: The heatmaps of clustering results (Left) and expression level of key genes (Right) obtained by the NMF, ONMF using DTPP and ONMF using proposed VS-ADMM applied to Dataset 2.