

Article

Tensor-Decomposition-Based Unsupervised Feature Extraction in Single-Cell Multiomics Data Analysis

Y-h. Taguchi ^{1,*}  and Turki Turki ² 

¹ Department of Physics, Chuo University, Tokyo 112-8551, Japan

² Department of Computer Science, King Abdulaziz University, Jeddah 21589, Saudi Arabia; tturki@kau.edu.sa

* Correspondence: tag@granular.com; Tel.: +81-3-3817-1791

Abstract: Analysis of single-cell multiomics datasets is a novel topic and is considerably challenging because such datasets contain a large number of features with numerous missing values. In this study, we implemented a recently proposed tensor-decomposition (TD)-based unsupervised feature extraction (FE) technique to address this difficult problem. The technique can successfully integrate single-cell multiomics data composed of gene expression, DNA methylation, and accessibility. Although the last two have large dimensions, as many as ten million, containing only a few percentage of nonzero values, TD-based unsupervised FE can integrate three omics datasets without filling in missing values. Together with UMAP, which is used frequently when embedding single-cell measurements into two-dimensional space, TD-based unsupervised FE can produce two-dimensional embedding coincident with classification when integrating single-cell omics datasets. Genes selected based on TD-based unsupervised FE are also significantly related to reasonable biological roles.

Keywords: tensor decomposition; feature extraction; single-cell; multiomics data



Citation: Taguchi, Y.-h.; Turki, T. Tensor-Decomposition-Based Unsupervised Feature Extraction in Single-Cell Multiomics Data Analysis. *Genes* **2021**, *12*, 1442. <https://doi.org/10.3390/genes12091442>

Academic Editors: Kenta Nakai and Tun-Wen Pai

Received: 25 August 2021

Accepted: 15 September 2021

Published: 18 September 2021

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

Single-cell multiomics data analysis is challenging [1]. There are multiple reasons for this issue. First, it inevitably includes too many missing values. In the usual high-throughput sequencing (HTS), the so-called depth can compensate for this problem. Nevertheless, because of the very limited amount of RNA retrieved from individual cells available, “depth” cannot resolve this missing value problem. Second, too many missing values result in apparent diversity. The primary purpose of single-cell analysis is to identify the diversity of individual cells that cannot be recognized by the tissue-level HTS. Although missing values are random, apparently very variant profiles appear from a single profile, which can be recognized if there is a large enough number of reads available. This compels researchers to distinguish between true biological diversity and apparent diversity caused by missing values [2].

Finally, single-cell analysis is computationally challenging. Because there are not many samples in the standard HTS, even if the number of features is large, the overall required computational resources decided by the product between the number of features and the number of samples are very limited. Nonetheless, since the number of samples that is the same as that of cells can be huge in single-cell analysis, single-cell analysis can be computationally very challenging.

To resolve these difficulties, we employed tensor-decomposition (TD)-based unsupervised feature extraction (FE) [3]. Prior to applying TD to multiomics datasets, singular-value decomposition (SVD) was applied to individual omics profiles such that individual omics profiles have common L singular-value vectors. Then, K omics profiles are formatted as an $L \times M \times K$ -dimensional tensor, where M is the number of single cells. Then, higher-order singular-value decomposition (HOSVD), which is a type of TD, is applied to the tensor. UMAP applied to singular-value vectors attributed to single cells by HOSVD success-

fully generated two-dimensional embedding, coincident with the known classification of single cells.

2. Materials and Methods

2.1. Gene Expression Profiles

Two single-cell multiomics datasets were downloaded from GEO using the following two GEO IDs.

2.1.1. GSE154762: Dataset 1

The multiomics dataset [4] retrieved from GEO ID GSE154762, which is denoted as Dataset 1 in this study, is composed of 899 single cells for which gene expression, DNA methylation, and DNA accessibility were measured. These single cells represent human oocyte maturation (Table 1). For gene expression, the file “GSE154762_hO_scChARM_count_matix.txt.gz” was downloaded from the supplementary file of GEO and was loaded into R [5] using the `read.table` function in R. For DNA methylation and DNA accessibility, 899 files with the extensions “WCG.bw” and “GCH.bw” were downloaded from the supplementary files of GEO and were loaded into R using the `import` function in the `rtracklayer` [6] package in R.

Table 1. The number of single cells within individual cell types included in Dataset 1.

FGO	GO1	GO2	Granulosa	Immune	MI	MII	StromaC1	StromaC2
81	40	46	93	20	155	90	189	185

2.1.2. GSE121708: Dataset 2

The multiomics dataset [7] retrieved from GEO ID GSE154762, which is denoted as Dataset 2 in this study, is composed of 852 single cells for which DNA methylation and DNA accessibility were measured, as well as 758 single cells for which gene expression was measured. These single cells represent the four time points of the mouse embryo (Table 2). For gene expression, the file “GSE121650_rna_counts.tsv.gz” was downloaded from the supplementary file of GEO and was loaded into R using the `read.table` function in R. For DNA methylation and DNA accessibility, 852 files with the extensions “met.tsv.gz” and “acc.tsv.gz” were downloaded from the supplementary file of GEO and were loaded into R using the `read.table` function in R.

Table 2. The number of single cells at four embryonic time points included in Dataset 2. For E7.5, the gene expression profiles of 296 single cells were measured.

E4.5-5.5	E6.5	E6.75	E7.5
267	98	97	390 (296)

2.2. Preprocessing of DNA Methylation Profiles

First, we collected genomic positions for which at least one measurement was performed for at least one single cell (i.e., union). Then, for each genomic position, three integers, -1 , 0 , and 1 , were assigned. When the genomic position was measured in a single cell and its state was methylated (nonmethylated), we attributed 1 (-1) to the genomic position of the single cell. Otherwise (i.e., missing observation), we attributed 0 to the genomic transition in a single cell. $x_{ij2} \in \mathbb{R}^{N_2 \times M}$ was stored as a sparse matrix object using the `Matrix` [8] package in R because of the large N_2 .

2.3. Preprocessing of DNA Accessibility

First, we divided the whole genome into 200 nucleotide regions, and DNA accessibility was summed up within individual regions. These values, which show the summation

of DNA accessibility within individual regions, are regarded as DNA accessibility at the individual 200 nucleotide regions, each of which is supposed to approximately correspond to a single nucleosome that is composed of 140-length DNA that wraps around histones and 80-length linker DNA. In this study, these 200 nucleotide regions are called “nucleosome regions”. $x_{ij3} \in \mathbb{R}^{N_3 \times M}$ was stored as a sparse matrix object using the Matrix package in R because of the large N_3 .

2.4. TD-Based Unsupervised FE

2.4.1. Reduction of Feature Dimensions

Here, feature denotes gene expression, DNA methylation, or DNA accessibility. Because the features of these three datasets differ from one another, we first applied SVD to these features. Suppose $x_{ijk} \in \mathbb{R}^{N_k \times M \times K}$ represents the value of the i th feature (expression of the i th gene, methylation of the i th genomic location, or DNA accessibility of the i th nucleosome region) at the k th single cell of the k th omics data ($1 \leq k \leq K = 3$, $k = 1$: gene expression, $k = 2$: DNA methylation, and $k = 3$: DNA accessibility). Applying SVD to x_{ijk} , we obtain:

$$x_{ijk} = \sum_{\ell=1}^L \lambda_{\ell} u_{\ell ik} v_{\ell jk} \quad (1)$$

where λ_{ℓ} is the ℓ th singular value and $u_{\ell ik}$ and $v_{\ell jk}$ are the i th and j th components of the ℓ th left and right singular-value vectors, respectively. Then, x_{ijk} is transformed to $x_{\ell jk} \in \mathbb{R}^{L \times M \times K}$ to have the same (common) feature dimension, L , independent of k , as:

$$x_{\ell jk} = \sum_{i=1}^{N_k} u_{\ell ik} x_{ijk} \quad (2)$$

2.4.2. Data Normalization

Prior to applying SVD to the individual omics profiles in these two datasets, x_{ijk} , ($k = 2, 3$), that is DNA methylation and accessibility, of Dataset 1 was normalized such that:

$$\sum_{i=1}^{N_k} |x_{ijk}| = N_k \quad (3)$$

whereas x_{ij1} , i.e., gene expression, was normalized such that:

$$\sum_{i=1}^{N_K} x_{ij1} = 0 \quad (4)$$

$$\sum_{i=1}^{N_k} x_{ij1}^2 = N_k \quad (5)$$

for Datasets 1 and 2. The reason why DNA methylation and the accessibility of Dataset 2 were not normalized is because $\sum_i |x_{ijk}|$, ($k = 2, 3$) is very small in some single cells in Dataset 2. Thus, applying normalization adds significant weight to these single cells with fewer observations and drastically skewed outcomes. To avoid this problem, x_{ijk} , ($k = 2, 3$) of Dataset 2 was not normalized.

2.4.3. TD Applied to Dimension-Reduced Multiomics Datasets

HOSVD [3] was applied to the tensor, $x_{\ell jk}$, and we obtained:

$$x_{\ell jk} = \sum_{\ell_1=1}^L \sum_{\ell_2=1}^M \sum_{\ell_3=1}^K G(\ell_1 \ell_2 \ell_3) u_{\ell_1 \ell} u_{\ell_2 j} u_{\ell_3 k} \quad (6)$$

where $G \in \mathbb{R}^{L \times M \times K}$ is the core tensor that represents the contribution of $u_{\ell_1 \ell} u_{\ell_2 j} u_{\ell_3 k}$ to $x_{\ell j k}$. $u_{\ell_1 \ell} \in \mathbb{R}^{L \times L}$, $u_{\ell_2 j} \in \mathbb{R}^{M \times M}$, $u_{\ell_3 k} \in \mathbb{R}^{K \times K}$ are singular-value matrices and are orthogonal matrices.

2.5. Categorical Regression

For categorical regression to test the coincidence between classification shown in Table 1 or Table 2 and singular-value vectors attributed to the j th single cells, we performed categorical regression:

$$v_{\ell j k} = a_{\ell k s} \delta_{j s} + b_{\ell k} \quad (7)$$

$$u_{\ell_2 j} = a_{\ell_2 s} \delta_{j s} + b_{\ell_2} \quad (8)$$

where s denotes one of the classifications shown in Table 1 or Table 2, $a_{\ell k s}$, $b_{\ell k}$, $a_{\ell_2 s}$, b_{ℓ_2} are regression coefficients, and $\delta_{j s}$ takes the value of 1 when the j th single cell belongs to the s th classification and 0 otherwise. Categorical regression was performed using the `ls` function in R. The obtained p -values were corrected using the Benjamini-Hochberg (BH) criterion [3]. ℓ_s or ℓ_{2s} associated with adjusted p -values less than 0.01 were regarded to be coincident with classification.

2.6. UMAP

Two-dimensional embedding was performed by UMAP [9]. The `umap` function implemented in R was used.

2.7. Gene Selection

After identifying which $u_{\ell_2 j}$ coincided with the classification, we needed to identify which $u_{\ell_1 \ell}$ was associated with the selected $u_{\ell_2 j}$ by investigating $|G(\ell_1 \ell_2 \ell_3)|$; ℓ_1 s with a larger $|G|$ with the selected ℓ_2 were regarded to be coincident with the classification. Then, the selected $u_{\ell_1 \ell}$ was converted back to $u_{\ell_1 i}$ attributed to genes as:

$$u_{\ell_1 i} = \sum_{\ell=1}^L u_{\ell_1 \ell} u_{\ell i} \quad (9)$$

p -values can be attributed to genes, i , assuming $u_{\ell_1 i}$ obeys a multiple Gaussian distribution (null hypothesis) as:

$$P_i = P_{\chi^2} \left[> \sum_{\ell_1} \left(\frac{u_{\ell_1 i}}{\sigma_{\ell_1}} \right)^2 \right] \quad (10)$$

where the summation is taken over only the selected ℓ_1 s, $P_{\chi^2}[> x]$ is the cumulative χ^2 distribution, where the argument is larger than x , and σ_{ℓ_1} is the standard deviation. P_i s were corrected by the BH criterion [3], and i s associated with adjusted P_i less than 0.01 were selected.

3. Results

3.1. Dataset 1

We obtained $x_{ij1} \in \mathbb{R}^{26500 \times 899}$, $x_{ij2} \in \mathbb{R}^{26438807 \times 899}$, and $x_{ij3} \in \mathbb{R}^{15478375 \times 899}$. SVD was applied to x_{ijk} with $L = 10$, as in Equation (1). For x_{ijk} , $k = 2, 3$, SVD was performed using the `irlba` function in the `irlba` package [10] in R because of the large N_k , $k = 2, 3$, as many as ten million. Then, HOSVD was applied to $x_{\ell j k}$, as in Equation (2).

One possible validation to check whether the above procedure works properly is to check whether $v_{\ell j k}$ and $u_{\ell_2 j}$ are coincident with the classification shown in Table 1. Because the above procedure is fully unsupervised, it is unlikely that $v_{\ell j k}$ and $u_{\ell_2 j}$ are accidentally coincident with the classification. To quantitatively validate the coincidence between the classification and $v_{\ell j k}$ or $u_{\ell_2 j}$, we applied categorical regression (see Section 2.5).

Table 3 shows the number of singular-value vectors coincident with the classification shown in Table 1. When SVD was applied to individual omics data, because $L = 10$, the number of singular-value vectors was 10 as well. When K omics datasets were integrated and HOSVD was applied, the number of singular-value vectors was $KL = 10K$. Thus, when DNA methylation and accessibility were integrated, 20 singular-value vectors were available. When all three omics data were integrated, 30 singular-value vectors were available. It is obvious that for all five cases, at least one singular-value vector was coincident with the classification. Thus, our strategy was essentially successful.

Table 3. Number of singular-value vectors coincident with classification shown in Table 1.

Adjusted <i>p</i> -Value	Gene Expression	SVD ($v_{\ell_{jk}}$)		HOSVD ($u_{\ell_{2j}}$)	
		DNA Methylation	DNA Accessibility	DNA Methylation and Accessibility	All
<0.01	10	7	1	10	18
≥0.01	0	3	9	10	12

To further validate the successful integration of singular-value vectors, we applied UMAP to 20 or 30 singular-value vectors obtained by HOSVD (Figure 1).

It is obvious that the integration of all three omics datasets (lower) was more coincident with classification than that of the integration of the two omics datasets, DNA methylation, and accessibility (upper). This suggests the usefulness of integrating the three omics datasets. In fact, single omics data cannot provide two-dimensional embedding coinciding with classification (Figure S1).

We also attempted to validate biological outcomes when all three omics datasets were integrated. We selected 47 genes associated with adjusted P_i less than 0.01, as described in Section 2.7 using u_{1i} because u_{1i} is associated with the largest:

$$\sum_{\ell_2} \sum_{\ell_3=1}^3 G^2(\ell_1 \ell_2 \ell_3), \quad (11)$$

where the summation of ℓ_2 is taken over only 18 ℓ_2 s coincident with the classification (Table 3). The selected 47 genes (Data S1) were uploaded to Enrichr [11].

Forty-seven genes were enriched by H3K36me3 based on “ENCODE Histone Modifications 2015”; H3K36m3 is known to play critical roles during oocyte maturation [12]. Forty-seven genes were also targeted by MYC based on “ENCODE and ChEA Consensus TFs from ChIP-X”; Myc is known to play critical roles in oogenesis [13]. Forty-seven genes were also targeted by TAF7 based on “ENCODE and ChEA Consensus TFs from ChIP-X” and “ENCODE TF ChIP-seq 2015”; TAF7 is known to play critical roles during oocyte growth [14]. Forty-seven genes were also targeted by ATF2 based on “ENCODE and ChEA Consensus TFs from ChIP-X”; the expression of ATF2 is known to be altered during oocyte development [15]. This suggests that our strategy correctly captures regulation-related parts (full data of the enrichment analysis are available as Data S1).

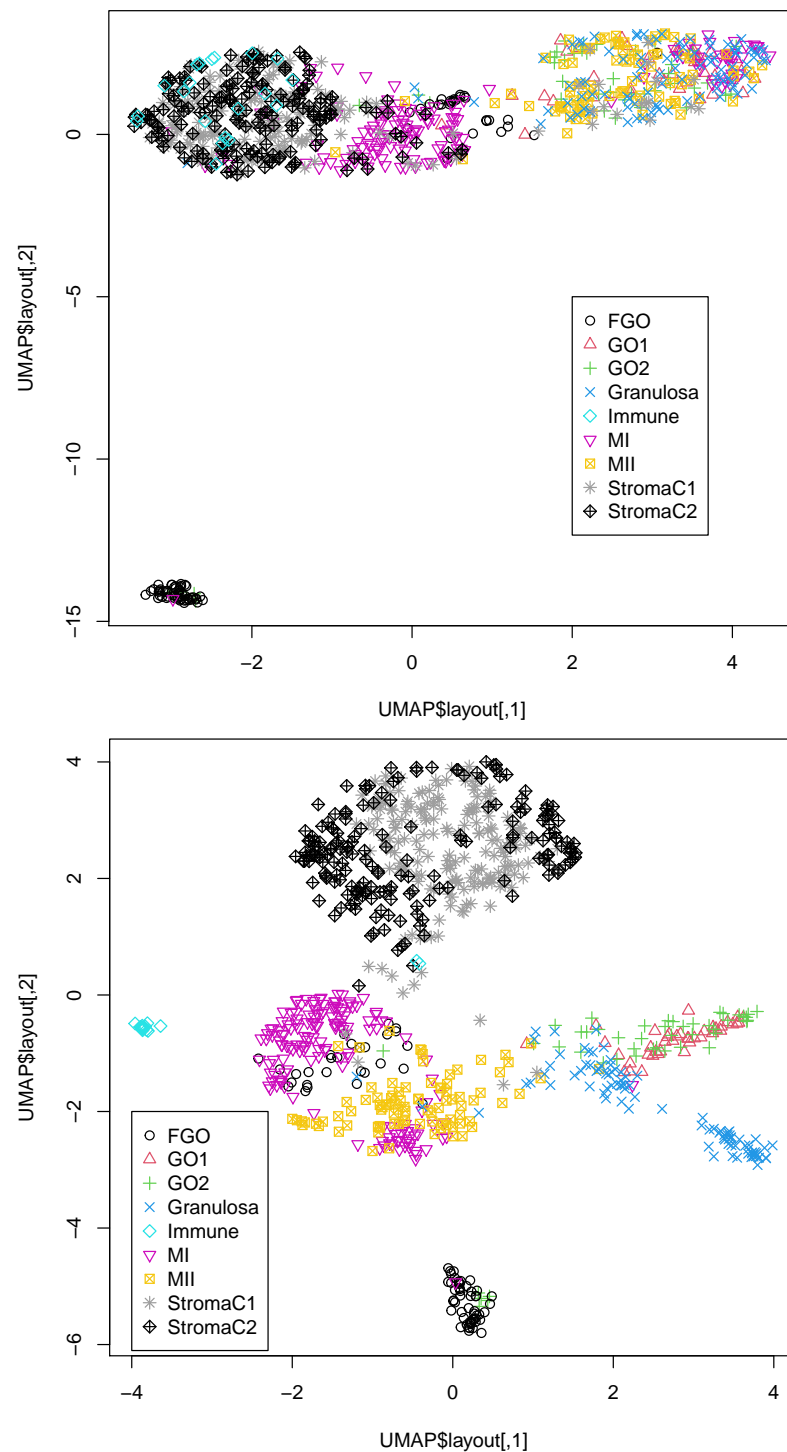


Figure 1. Two-dimensional embedding of singular-value vectors, $u_{\ell_2 j}$, computed by HOSVD applied to $x_{\ell j k}$ in Dataset 1 (Table 3). Upper: $u_{\ell_2 j}, 1 \leq \ell_2 \leq 20$ when only DNA methylation and accessibility ($k = 2, 3$) are integrated. Lower: $u_{\ell_2 j}, 1 \leq \ell_2 \leq 30$ when all three omics data points ($1 \leq k \leq 3$) are integrated. Default settings other than `custom.config$n_neighbors = 100` were used.

3.2. Dataset 2

To confirm that the success in the previous section was not accidental, we applied the same procedure to Dataset 2 as well. We obtained $x_{ij1} \in \mathbb{R}^{22084 \times 758}$, $x_{ij2} \in \mathbb{R}^{20106507 \times 852}$, and $x_{ij3} \in \mathbb{R}^{13627678 \times 852}$. SVD was applied to x_{ijk} with $L = 10$, as in Equation (1). For $x_{ijk}, k = 2, 3$, SVD was performed using the `irlba` function in the `irlba` package [10] in R because of the large $N_k, k = 2, 3$ of as many as ten million. Then, HOSVD was applied to

$x_{\ell_{jk}}$, as in Equation (2). Because $N_1 = 758 < N_2 = N_3 = 852$, when HOSVD was applied to $x_{\ell_{jk}}$ composed of all three omics datasets, only 758 single cells shared with all x_{ijk} were considered. As described in the previous section, we first validated the coincidence between singular-value vectors attributed to single cells (Table 4), that is v_{ℓ_j} and $u_{\ell_2\ell_k}$, and the classification in Table 2.

Table 4. Number of singular-value vectors coincident with the classification shown in Table 2.

Adjusted <i>p</i> -Value	Gene Expression	SVD ($v_{\ell_{jk}}$)		HOSVD (u_{ℓ_2j})	
		DNA Methylation	DNA Accessibility	DNA Methylation and Accessibility	All
<0.01	10	7	5	10	18
≥0.01	0	3	5	10	12

The coincidence between the singular-value vectors and the classification in Table 4 was even better than that in Table 3. Thus, it is unlikely that the superior outcome in Table 3 was purely accidental. To further validate the successful integration of singular-value vectors, we applied UMAP to 20 or 30 singular-value vectors obtained by HOSVD (Figure 2).

It is obvious that the integration of all three omics datasets (lower) was more coincident with classification than that of the integration of the two omics datasets, DNA methylation, and DNA accessibility (upper), as can be seen in Figure 2. This again confirms the usefulness of integrating the three omics datasets. In fact, single omics data cannot provide two-dimensional embedding coinciding with classification (Figure S2).

We also attempted to validate biological outcomes when all three omics datasets were integrated. We selected 175 genes associated with adjusted P_i less than 0.01, as described in Section 2.7 using u_{1i} because u_{1i} is associated with the largest:

$$\sum_{\ell_2} \sum_{\ell_3=1}^3 G^2(\ell_1\ell_2\ell_3) \quad (12)$$

where the summation of ℓ_2 is taken over only 18 ℓ_2 s coincident with the classification (Table 4). The selected 175 genes (Data S2) were converted to gene symbols by the DAVID [16] gene ID converter and were uploaded to Enrichr.

One-hundred and seventy-five genes were enriched by H3K36me3 based on “ENCODE Histone Modifications 2015”; H3K36m3 is known to play critical roles during gastrulation [17]. One-hundred and seventy-five genes were also targeted by MYC based on “ENCODE and ChEA Consensus TFs from ChIP-X”; Myc is also known to play critical roles in gastrulation [18]. One hundred and seventy-five genes were also targeted by TAF7 based on “ENCODE and ChEA Consensus TFs from ChIP-X” and “ENCODE TF ChIP-seq 2015”; TAF7 is known to play critical roles during gastrulation [19]. One-hundred and seventy-five genes were also targeted by ATF2 based on “ENCODE and ChEA Consensus TFs from ChIP-X”; the expression of ATF2 is known to be maintained during gastrulation [20]. This suggests that our strategy correctly captured regulation-related parts (full data of the enrichment analysis are available as Data S2).

These two examples, the application to Datasets 1 and 2, demonstrate the usefulness of the present strategy to integrate single-cell multiomics datasets composed of gene expression, DNA methylation, and accessibility.

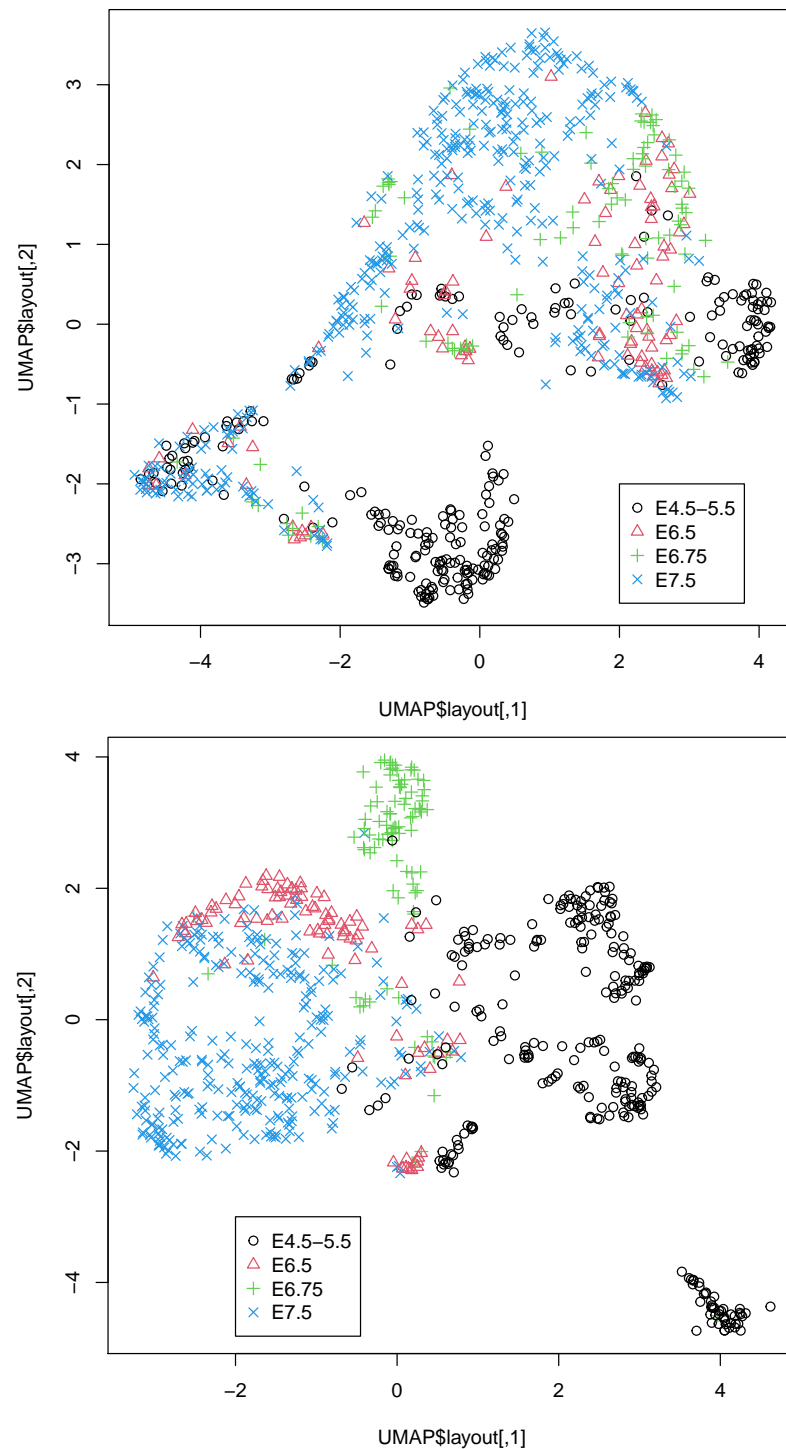


Figure 2. Two-dimensional embedding of singular-value vectors, $u_{\ell_2 j}$, computed by HOSVD applied to $x_{\ell j k}$ in Dataset 2 (Table 4). Upper: $u_{\ell_2 j}, 1 \leq \ell_2 \leq 20$ when only DNA methylation and accessibility ($k = 2, 3$) are integrated. Lower: $u_{\ell_2 j}, 1 \leq \ell_2 \leq 30$ when all three omics data points ($1 \leq k \leq 3$) are integrated. Default settings other than `custom.config$n_neighbors = 100` were used.

4. Discussion

In this study, we demonstrated the usefulness of our strategy when it was applied to the integrated analysis of single-cell multiomics datasets composed of gene expression, DNA methylation, and DNA accessibility. One might wonder if other more popular methods can achieve similar performance because our strategy is useless if others can

perform comparably. There are several advantages of our method, which other methods do not have.

First, we do not have to fill in the missing values with nonzero values. Single-cell measurements are usually associated with a large number of missing values (Table 5).

Table 5. Number of single cells, features, nonzero components, and their ratios.

Numbers	Expression	DNA Methylation	DNA Accessibility
Dataset 1			
single cells	899	899	899
features	26,500	26,438,807	15,478,375
total components	2.38×10^7	2.38×10^{10}	1.39×10^{10}
nonzero components	6.76×10^6	5.50×10^8	3.85×10^8
the ratio of nonzero components	0.28	0.02	0.03
Data set 2			
single cells	758	852	852
features	22,084	20,106,507	13,627,678
total components	1.67×10^7	1.71×10^{10}	1.16×10^{10}
nonzero components	4.87×10^6	6.96×10^8	7.87×10^8
the ratio of nonzero components	0.29	0.04	0.07

Although gene expression profiles were associated with a relatively small number of missing components, more than 70 % were missing. For DNA methylation and accessibility, the situation was very difficult to treat. Only a few percentages of components had values, while the rest were missing values. To address this problem, especially for DNA methylation and accessibility, heavy preprocessing is usually required. For example, for Dataset 1, statistical tests were applied and regions associated with significant p -values were selected [4], which reduced the number of features attributed to DNA methylation and accessibility. Because such a statistical test automatically filters out regions filled in with missing values, the ratio of nonzero components was also reduced as a result. For Dataset 2, the authors restricted the features to only the most variable ones (typically $\sim 10^3$) and occasionally filled in missing components with Bayesian models [7]. These procedures inevitably introduce arbitrariness to the outcomes, as preprocessing the data might affect the outcome. In contrast to these arbitrary procedures, our method is almost unsupervised. We did not select any features or fill in the missing values. Despite these fully unsupervised strategies, our results were highly coincident with the classification (Tables 3 and 4 and Figures 1 and 2). From this perspective, our strategy is superior to the other methods.

Second, our method can deal with massive datasets. For example, although integrated analysis of multiomics data was performed using multiomics factor analysis (MOFA) [21] in the original studies [4,7] of Datasets 1 and 2, MOFA cannot accept x_{ijk} in this study as inputs because MOFA does not implement sparse matrix architecture. During the computation of MOFA, zero values must be filled in with nonzero values to evaluate the convergence; this results in a dense matrix that cannot be stored in the computer memory because the number of components of DNA methylation and accessibility is too large to store them as they are (Table 5). In our computation, we can apply SVD to these large datasets while keeping them in a sparse matrix format using the `irlba` package implemented in R. SVD not only reduces the number of features to L , but also fills in missing values. Thus, we can manage a large matrix as in our implementation.

Third, our method is free from the dividing weight between multiomics datasets; how to weigh individual omics data must be decided based on some criteria outside the datasets available. Nevertheless, in our implementation, the weight of individual omics data is represented by $u_{\ell_3 k}$, which is automatically decided by simply applying HOSVD to a multiomics dataset. Thus, from this perspective, our strategy is outstanding.

Although we showed that the integration of all three omics data was superior to that of the integration of DNA methylation and accessibility (Figures 1 and 2), one might wonder if the integration of gene expression and DNA methylation or DNA accessibility might be comparable to that of all three omics datasets. In order to deny this possibility, we also considered these combinations of two of the three omics datasets (Figures S3 and S4). Although the integration of gene expression and DNA accessibility in Dataset 1 (Figure S3B) is comparable to that of all three omics data, neither integration of gene expression and DNA methylation (Figure S4A) nor that of gene expression and DNA accessibility (Figure S4B) is comparable to that of all three omics data in Dataset 2. Thus, it is obvious that only the integration of the three omics datasets can give us UMAP embedding coincident with the classification regardless of the dataset considered.

As for the comparisons with other methods, as mentioned above, no methods implemented with a sparse matrix architecture and applicable to multiomics datasets exist to our knowledge. Thus, we could not compare our performance to other methods.

Prospective uses of our methods are as follows. First of all, it can integrate gene expression profiles, DNA methylation, and accessibility in single-cell measurements without applying preprocessing; this enables researchers to obtain reasonable results without struggling to convert raw data into treatable formats. In addition to this, since it can save the memories required for analyzing single-cell multiomics datasets, more researchers who do not have massive computational facilities can analyze massive single-cell measurements.

5. Conclusions

In this study, we proposed a method for applying TD-based unsupervised FE to single-cell multiomics datasets composed of gene expression, DNA methylation, and DNA accessibility. Together with UMAP, the proposed method successfully integrated a multiomics dataset and generated a two-dimensional embedding of single cells coincident with the classification. The implementation requires neither filling missing values nor massive CPU memory to store multiomics datasets of single cells and can deal with DNA methylation and accessibility with ten million features. The present implementation is very promising and can be a de facto standard method to integrate single-cell multiomics datasets composed of gene expression, DNA methylation, and DNA accessibility.

Supplementary Materials: The following are available online at <https://www.mdpi.com/article/10.3390/genes12091442/s1>, Figure S1: UMAP embedding of single omics data for data set 1, Figure S2: UMAP embedding of single omics data for data set 2, Figure S3: UMAP embedding of integration of two omics data for data set 1, Figure S4: UMAP embedding of integration of two omics data for data set 2.

Author Contributions: Y.-h.T. planned the research and performed the analyses. Y.-h.T. and T.T. evaluated the results, discussions, and outcomes and wrote and reviewed the manuscript. All authors have read and agreed to the published version of the manuscript.

Funding: This work was supported by KAKENHI (Grant Numbers 19H05270, 20H04848, and 20K12067) to Y.-h.T.

Data Availability Statement: The data used in this study are available in GEO ID GSE154762 and GSE121708. A sample of the R source can be found at <https://github.com/tagtag/scMultiR> (accessed on 14 September 2021).

Conflicts of Interest: The authors declare no conflict of interest. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript; nor in the decision to publish the results.

Abbreviations

The following abbreviations are used in this manuscript:

BH	Benjamini–Hochberg
FE	feature extraction
HOSVD	higher-order singular-value decomposition
HTS	high-throughput sequencing
MOFA	Multi-Omics Factor208Analysis
SVD	singular-value decomposition
TD	tensor decomposition

References

- Lee, J.; Hyeon, D.Y.; Hwang, D. Single-cell multiomics: Technologies and data analysis methods. *Exp. Mol. Med.* **2020**, *52*, 1428–1442. [[CrossRef](#)] [[PubMed](#)]
- Liu, Q.; Herring, C.A.; Sheng, Q.; Ping, J.; Simmons, A.J.; Chen, B.; Banerjee, A.; Li, W.; Gu, G.; Coffey, R.J.; et al. Quantitative assessment of cell population diversity in single-cell landscapes. *PLoS Biol.* **2018**, *16*, 1–29. [[CrossRef](#)] [[PubMed](#)]
- Taguchi, Y.H. *Unsupervised Feature Extraction Applied to Bioinformatics*; Springer International Publishing: Berlin/Heidelberg, Germany, 2020; [[CrossRef](#)]
- Yan, R.; Gu, C.; You, D.; Huang, Z.; Qian, J.; Yang, Q.; Cheng, X.; Zhang, L.; Wang, H.; Wang, P.; et al. Decoding dynamic epigenetic landscapes in human oocytes using single-cell multi-omics sequencing. *Cell Stem Cell* **2021**, *28*, 1641–1656.e7. [[CrossRef](#)] [[PubMed](#)]
- R Core Team. *R: A Language and Environment for Statistical Computing*; R Foundation for Statistical Computing: Vienna, Austria, 2020.
- Lawrence, M.; Gentleman, R.; Carey, V. rtracklayer: An R package for interfacing with genome browsers. *Bioinformatics* **2009**, *25*, 1841–1842. [[CrossRef](#)] [[PubMed](#)]
- Argelaguet, R.; Clark, S.J.; Mohammed, H.; Stapel, L.C.; Krueger, C.; Kapourani, C.A.; Imaz-Rosshandler, I.; Lohoff, T.; Xiang, Y.; Hanna, C.W.; et al. Multi-omics profiling of mouse gastrulation at single-cell resolution. *Nature* **2019**, *576*, 487–491. [[CrossRef](#)] [[PubMed](#)]
- Bates, D.; Maechler, M. *Matrix: Sparse and Dense Matrix Classes and Methods*; R Package Version 1.3-4; R Package: Boston, MA, USA, 2021.
- Becht, E.; McInnes, L.; Healy, J.; Dutertre, C.A.; Kwok, I.W.H.; Ng, L.G.; Ginhoux, F.; Newell, E.W. Dimensionality reduction for visualizing single-cell data using UMAP. *Nat. Biotechnol.* **2018**, *37*, 38–44. [[CrossRef](#)] [[PubMed](#)]
- Baglama, J.; Reichel, L.; Lewis, B.W. *irlba: Fast Truncated Singular Value Decomposition and Principal Components Analysis for Large Dense and Sparse Matrices*; R Package Version 2.3.3; R Package: Boston, MA, USA, 2019.
- Kuleshov, M.V.; Jones, M.R.; Rouillard, A.D.; Fernandez, N.F.; Duan, Q.; Wang, Z.; Koplev, S.; Jenkins, S.L.; Jagodnik, K.M.; Lachmann, A.; et al. Enrichr: A comprehensive gene set enrichment analysis web server 2016 update. *Nucleic Acids Res.* **2016**, *44*, W90–W97. [[CrossRef](#)] [[PubMed](#)]
- Xu, Q.; Xiang, Y.; Wang, Q.; Wang, L.; Brind'Amour, J.; Bogutz, A.B.; Zhang, Y.; Zhang, B.; Yu, G.; Xia, W.; et al. SETD2 regulates the maternal epigenome, genomic imprinting and embryonic development. *Nat. Genet.* **2019**, *51*, 844–856. [[CrossRef](#)] [[PubMed](#)]
- Suzuki, T.; Ichiro Abe, K.; Inoue, A.; Aoki, F. Expression of c-MYC in Nuclear Speckles During Mouse Oocyte Growth and Preimplantation Development. *J. Reprod. Dev.* **2009**, *55*, 491–495. [[CrossRef](#)] [[PubMed](#)]
- Yu, C.; Cvetesic, N.; Gupta, K.; Ye, T.; Gazdag, E.; Hisler, V.; Negroni, L.; Hajkova, P.; Lenhard, B.; Müller, F.; et al. TBPL2/TFIIA complex overhauls oocyte transcriptome during oocyte growth. *bioRxiv* **2020**. [[CrossRef](#)]
- Vigneault, C.; McGraw, S.; Sirard, M.A. Spatiotemporal expression of transcriptional regulators in concert with the maternal-to-embryonic transition during bovine in vitro embryogenesis. *Reproduction* **2009**, *137*, 13–21. [[CrossRef](#)] [[PubMed](#)]
- Huang, D.W.; Sherman, B.T.; Lempicki, R.A. Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. *Nat. Protoc.* **2008**, *4*, 44–57. [[CrossRef](#)] [[PubMed](#)]
- Sitbon, D.; Boyarchuk, E.; Dingli, F.; Loew, D.; Almouzni, G. Histone variant H3.3 residue S31 is essential for *Xenopus* gastrulation regardless of the deposition pathway. *Nat. Commun.* **2020**, *11*, 1256. [[CrossRef](#)] [[PubMed](#)]
- Downs, K.M.; Martin, G.R.; Bishop, J.M. Contrasting patterns of myc and N-myc expression during gastrulation of the mouse embryo. *Genes Dev.* **1989**, *3*, 860–869. [[CrossRef](#)] [[PubMed](#)]
- Langer, D.; Martianov, I.; Alpern, D.; Rhinn, M.; Keime, C.; Dollé, P.; Mengus, G.; Davidson, I. Essential role of the TFIID subunit TAF4 in murine embryogenesis and embryonic stem cell differentiation. *Nat. Commun.* **2016**, *7*, 11063. [[CrossRef](#)] [[PubMed](#)]
- Villarreal, X.C.; Richter, J.D. Analysis of ATF2 gene expression during early *Xenopus laevis* development. *Gene* **1995**, *153*, 225–229. [[CrossRef](#)]
- Argelaguet, R.; Velten, B.; Arnol, D.; Dietrich, S.; Zenz, T.; Marioni, J.C.; Buettner, F.; Huber, W.; Stegle, O. Multi-Omics Factor Analysis—A framework for unsupervised integration of multi-omics datasets. *Mol. Syst. Biol.* **2018**, *14*, e8124. [[CrossRef](#)] [[PubMed](#)]