

RESEARCH

Open Access

ScLRTC: imputation for single-cell RNA-seq data via low-rank tensor completion



Xiutao Pan, Zhong Li*, Shengwei Qin, Minzhe Yu and Hang Hu

Abstract

Background: With single-cell RNA sequencing (scRNA-seq) methods, gene expression patterns at the single-cell resolution can be revealed. But as impacted by current technical defects, dropout events in scRNA-seq lead to missing data and noise in the gene-cell expression matrix and adversely affect downstream analyses. Accordingly, the true gene expression level should be recovered before the downstream analysis is carried out.

Results: In this paper, a novel low-rank tensor completion-based method, termed as scLRTC, is proposed to impute the dropout entries of a given scRNA-seq expression. It initially exploits the similarity of single cells to build a third-order low-rank tensor and employs the tensor decomposition to denoise the data. Subsequently, it reconstructs the cell expression by adopting the low-rank tensor completion algorithm, which can restore the gene-to-gene and cell-to-cell correlations. ScLRTC is compared with other state-of-the-art methods on simulated datasets and real scRNA-seq datasets with different data sizes. Specific to simulated datasets, scLRTC outperforms other methods in imputing the dropouts closest to the original expression values, which is assessed by both the sum of squared error (SSE) and Pearson correlation coefficient (PCC). In terms of real datasets, scLRTC achieves the most accurate cell classification results in spite of the choice of different clustering methods (e.g., SC3 or t-SNE followed by K-means), which is evaluated by using adjusted rand index (ARI) and normalized mutual information (NMI). Lastly, scLRTC is demonstrated to be also effective in cell visualization and in inferring cell lineage trajectories.

Conclusions: a novel low-rank tensor completion-based method scLRTC gave imputation results better than the state-of-the-art tools. Source code of scLRTC can be accessed at <https://github.com/jianghuaijie/scLRTC>.

Keywords: Single-cell RNA-seq, Data imputation, Low-rank tensor

Background

Over the past few years, with the explosive growth of scRNA sequence data, important biological discoveries have been progressively conducted. However, as impacted by the picogram level of RNAs in a single cell, RNA transcripts may be missed during the reverse transcription and amplification step, so the transcripts are not detected in the following sequencing, which is termed as the dropout problem [1]. The resulting gene-cell expression matrix will consist of numerous false zeros attributed to dropout events, which will corrupt

the biological signal and impede downstream analyses (e.g., cell clustering, data visualization and cell trajectory inference). To reduce the impact of this problem, besides increasing the efficiency of transcription capture, an effective imputation algorithm for scRNA-seq data should be developed to predict missing values attributed to dropout events [2].

Existing single-cell imputation methods have two main types: one complies with the deep learning method. For instance, DeepImpute [3] was designed to impute the scRNA sequence by applying a deep neural network (DNN) with a dropout layer and loss function to learn patterns in the data. DCA [4] established an auto-encoder to model the distribution of genes with a zero-

* Correspondence: lizhong@zstu.edu.cn

Department of Mathematical Sciences, School of Science, Zhejiang Sci-Tech University, Hangzhou 310018, China



© The Author(s). 2021 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

inflated negative binomial prior, and then attempted to predict the mean, standard deviation and dropout probability of genes. ScIGANs [5] adopted the generative adversarial networks (GANs) to learn the dependence of nonlinear genes and genes from complex multi-cell type samples, and then trained the neural network model to generate real expression profiles of defined cell types. However, due to the influence of the training set and the existence of over-fitting problem, these methods may generate the false-positive results in differential expression analyses [6]. Another type for single-cell imputation methods complies with the statistical algorithm. For instance, SAVER [7] exploited information across genes in the identical cell type with a Bayesian approach to recover true expression levels; it also measured the uncertainty of recovered values. MAGIC [8] performed a soft clustering after building a Markov transition matrix, and then replaced a gene's raw expression with its weighted mean expression in a cluster. However, MAGIC also imputes the gene expression values that are not affected by dropout. Therefore, it may introduce the new bias into the data and possibly eliminate the meaningful biological variations. ScImpute [9] initially estimated the probability of an entry to be dropout with the use of a mixture model, and then imputed the potential dropout entries of a cell by employing the information from the gene expression of consistent cells. DrImpute [10] presented a clustering-based method and implemented a consensus strategy which estimated a value with several cluster priors or distance matrices and then imputed the data by aggregation. CMF-Impute [11] drew upon the similarity of cells and genes to build a collaborative matrix factorization-based model for imputing the dropout entries of a given scRNA-seq expression. ALAR [12] provided a low-rank approximation of the expression matrix using singular vector decomposition (SVD). McImpute [13] used the nuclear norm minimization to realize a matrix completion algorithm for the scRNA data imputation. A study [14] suggested that taking advantage of the presence of low-rank submatrix can improve the imputation performance compared to the traditional low-rank matrix restore methods. For example, PBLR [15] considered the cell grouping information and performed a bounded low-rank completion method for each group. ScLRTD [16] introduced the tensor into the imputation of single-cell datasets, but it is mainly aimed at the completion with single-cell multi-omics sequencing data and the result in the scRNA dataset is not better than MAGIC, because this tensor based method did not fully take advantage of the correlation of single-cell data. Liu et al. [17] proposed a definition for the tensor trace norm that generalizes the established definition of matrix trace norm. Similar to the matrix based imputation, the tensor based imputation is

formulated as a convex optimization problem and is solved by three algorithms SiLRTC, FaLRTC and HaLRTC. Experimental comparisons show that these methods are more accurate and robust than other heuristic approaches (Tucker, Parafac and SVD), which can propagate the data structure to fill large missing regions.

In this paper, a novel low-rank tensor completion based method (scLRTC) is proposed for the scRNA-seq data imputation. Since scRNA-seq data commonly involve single cells from different cell types and single cells with the identical type normally exhibit the similar expression pattern, the underlying true expression matrix is reasonably assumed to be able to be approximated by a low-rank matrix [1]. Based on such an assumption, the similar expression patterns of a single cell are adopted to build a third-order tensor, and then the single-cell gene expression is restored by approximating the tensor rank. This method is applied to nine scRNA-seq datasets and four simulation datasets, and it is compared with several state-of-the-art methods (SAVER [7], MAGIC [8], scImpute [9], DrImpute [10], CMF-Impute [11], PBLR [15], WEDGE [18] and scGNN [19]). As revealed from considerable data analyses, the proposed method is capable of achieving more accurate imputation results and improving the downstream analysis.

Results

The proposed imputation method is employed for nine published scRNA-seq datasets (i.e., Pollen [20], Usoskin [21], Yan [22], Zeisel [23], Mouse [24], PBMC [25], Chen [26], Loh [27] and Petropoulos [28]) and four simulation datasets generated from Splatter package [29], and it is compared with some popular methods (e.g., SAVER [7], MAGIC [8], scImpute [9], DrImpute [10], CMF-Impute [11], PBLR [15], WEDGE [18] and scGNN [19]). The method is largely evaluated from five aspects with cell subpopulation clustering, dimensionality reduction visualization, data masking evaluation, correlation analysis and differential expression analysis, and cell trajectory inference. The parameter setting of the respective dataset is listed in Table 1, and the parameter setting of the simulation dataset is presented in Table 2.

Evaluating imputation accuracy through cell clustering

In the relevant research on the scRNA-seq dataset, cell clustering refers to one of the critical contents. There are many clustering algorithms (e.g., K-means and SC3 [30]). Among the mentioned methods, SC3 is recognized as an accurate unsupervised single-cell clustering tool that does not explicitly address dropout events for the scRNA-seq data. Thus, the proposed method together with other popular scRNA-seq imputation methods was added into the preprocessing step of SC3. Then we used the cell clustering accuracy measured by adjusted rand

Table 1 Parameter settings corresponding to respective datasets in our experiment

Dataset	K	P	α	ρ	epsilon
Pollen	10	10	[1 1e-1 2e-3]	1e-5	1e-2
Usoskin	8	8	[1 1e-2 1e-3]	1e-5	1e-2
Yan	5	5	[1 1e-2 2e-3]	1e-5	1e-2
Zeisel	20	20	[1 1e-2 1e-3]	1e-5	1e-2
Mouse	10	10	[1 1e-2 1e-3]	1e-5	1e-2
PBMC	10	10	[1 1e-2 1e-3]	1e-5	1e-3
Chen	10	10	[1 1e-2 1e-3]	1e-5	1e-3
Loh	11	11	[1 1e-2 2e-3]	2e-4	1e-2
Petropoulos	10	10	[1 1e-2 1e-3]	1e-5	1e-2
Simulation dataset	15	8	[1 1e-1 1e-3]	1e-5	1e-5

index (ARI) [31] and normalized mutual information (NMI) [32] to evaluate their performance, namely, the consistency between the inferred cell cluster and the real cell cluster.

After the data imputation by the proposed scLRTC and other methods (DrImpute, SAVER, scImpute, MAGIC, CMF-Impute and PBLR), we used SC3 to cluster 6 published scRNA-seq datasets, including Usoskin, Pollen, Yan, Zeisel, Mouse and PBMC. The clustering accuracy measured by ARI and NMI are plotted in Fig. 1A and Fig. 1B, respectively. Obviously, the proposed method has the best ARI performance in Usoskin, Pollen, Yan, Zeiel, Mouse and PBMC, and the performance of NMI on the Usoskin dataset can be as competitive as CMF-Impute. In summary, the proposed method imputation can improve the clustering accuracy of SC3.

To show that the proposed imputation method does not depend on the clustering method, we further used another popular single-cell clustering method (first using the dimensionality reduction by t-distributed stochastic neighbor embedding (t-SNE) [33], and then applying K-means for clustering) [10] to test the performance of the proposed scLRTC and other methods. Compared with the SC3 algorithm, the K-means algorithm is more affected by the initial values. To compare the clustering results more reasonably, we performed t-SNE + K-means 20 times on the Pollen and Usoskin datasets (The

perplexity of t-SNE is set to 10, and other parameters of t-SNE are set as default parameters). In the Pollen dataset, our median is the highest at 0.722, and the maximum value is 0.853, which is better than the maximum value (0.847) by DrImpute (Fig. 2A). In the Usoskin dataset, the median of the proposed method is 0.684 and the maximum value is 0.742, which are both the highest compared to other methods (Fig. 2B). Note that the result of scImpute is worse than that of SAVER. When verifying the clustering performance by SC3 and t-SNE + K-means clustering, similar result was also appeared in CMF-Impute [11]. The main reason may be that scImpute relies on the spectral clustering which may influence the subsequent imputation process when the data has unbalanced clusters. In brief, the proposed imputation achieves a better overall effect than other imputation algorithms.

Furthermore, we compared scLRTC with the latest matrix completion based method WEDGE [18] and deep learning based method scGNN [19] on the Zeisel dataset. We applied the Scanpy's Louvain algorithm [34, 35] for the scRNA-seq data clustering and found scLRTC achieved an ARI of 0.692, which is higher than WEDGE's 0.560 and scGNN's 0.678. Finally, we did the test for a large Chen dataset [26] where the number of cells is more than 10,000. We also used the Louvain algorithm to cluster the scRNA-seq data and found the clustering performance index ARI increased from 0.611 (raw data) to 0.673 by scLRTC. Considering the time complexity of scLRTC, the addition of tensor computation makes it slower than other methods. But we can control the size of tensor for various datasets to relief the influence of tensor computation. Figure 2C illustrates the running time of scLRTC for the mentioned experimental datasets with different sizes of tensor setting. It shows that the time complexity of scLRTC is not quadratic proportional to the number of cells, which makes it applicable for scRNA-seq datasets with different sizes.

Cell visualization

Visually representing scRNA-seq data involves shrinking the gene expression matrix into a lower space, and then

Table 2 Parameter settings of simulation datasets generated from Splatter

Parameter	Simulation dataset	Parameter	Simulation dataset
version	1.10.1	dropout.type	"group"
nGenes	1000	method	"groups"
nCells	500	de.prob	c(0.05, 0.08, 0.01)
group.prob	c(0.3, 0.3, 0.4)	de.facLoc	0.5
dropout.shape	c(ds, ds, ds), ds ∈ {-0.3, 0, 0.05, 0.25}	de.facScale	0.8
dropout.mid	Default	dropout.present	Null

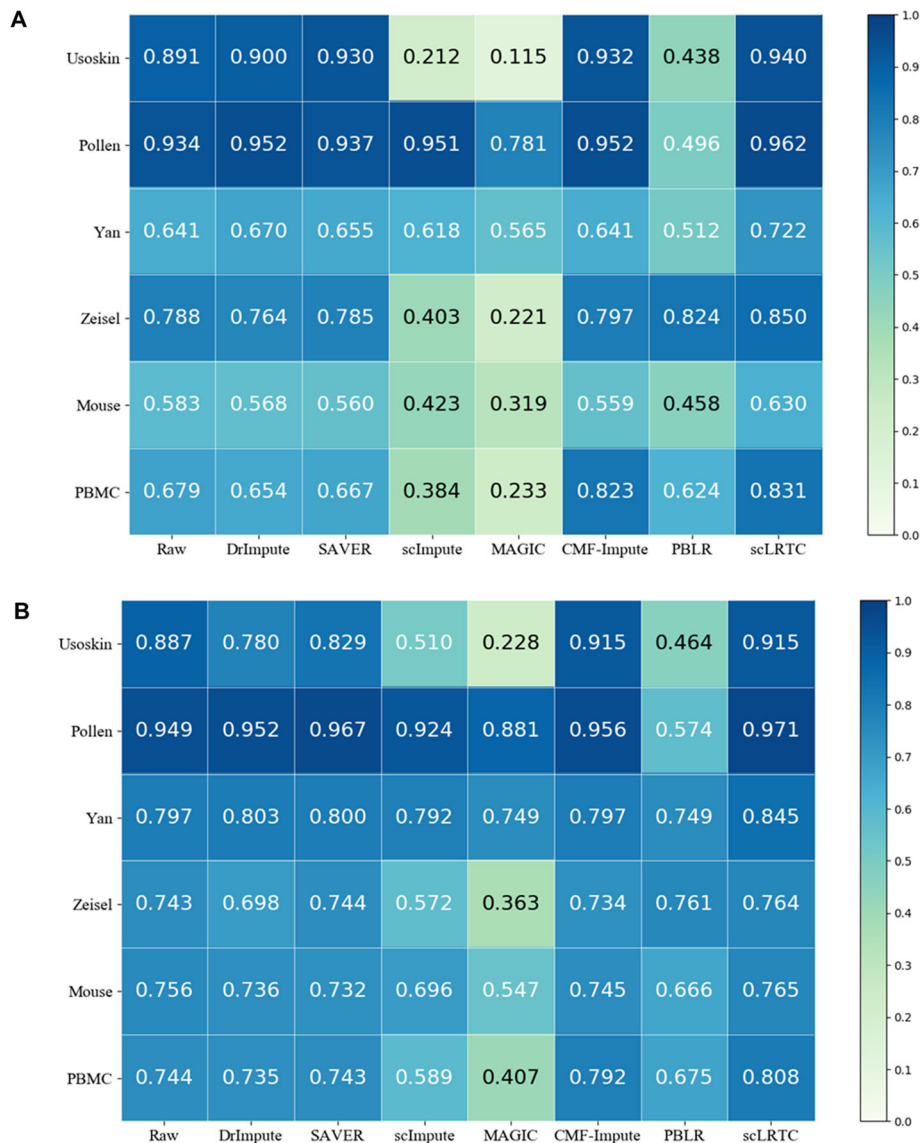


Fig. 1 SC3 clustering comparison of scLRTC and other methods on different datasets. **(A)** Comparison of ARI indicators obtained by SC3 clustering on 6 scRNA-seq datasets using different algorithms. **(B)** Comparison of NMI indicators obtained by SC3 clustering on 6 scRNA-seq datasets using different algorithms

mapping each cell’s transcriptome in the reduced low dimensional space. Several dimensionality reduction methods are generally known (e.g., PCA [36], t-SNE and UMAP [37]), where UMAP is suggested to be particularly suitable for the visualization of any dimensional data. Accordingly, UMAP was employed to discuss the dimensionality reduction effect before and after imputation on four expression matrices of Yan, Pollen, Usoskin and Zeisel datasets. To be specific, cells were visualized in a two-dimensional space, and different cells were stained using real labels before and after imputation. To quantify the grouping of cell

transcriptomes, an unsupervised clustering quality measurement was conducted with silhouette coefficient (SC) [38] to evaluate the effect of dimensionality reduction. The higher the silhouette coefficient, the more significant the dimensionality reduction effect will be. The UMAP dimensionality reduction visualization and the average SC of the raw and imputed data (4 published datasets) with different methods are illustrated in Fig. 3 and Fig. 4. According to these figures, the SC values of the proposed scLRTC in these datasets are the highest with 0.884, 0.797, 0.861, 0.639, respectively.

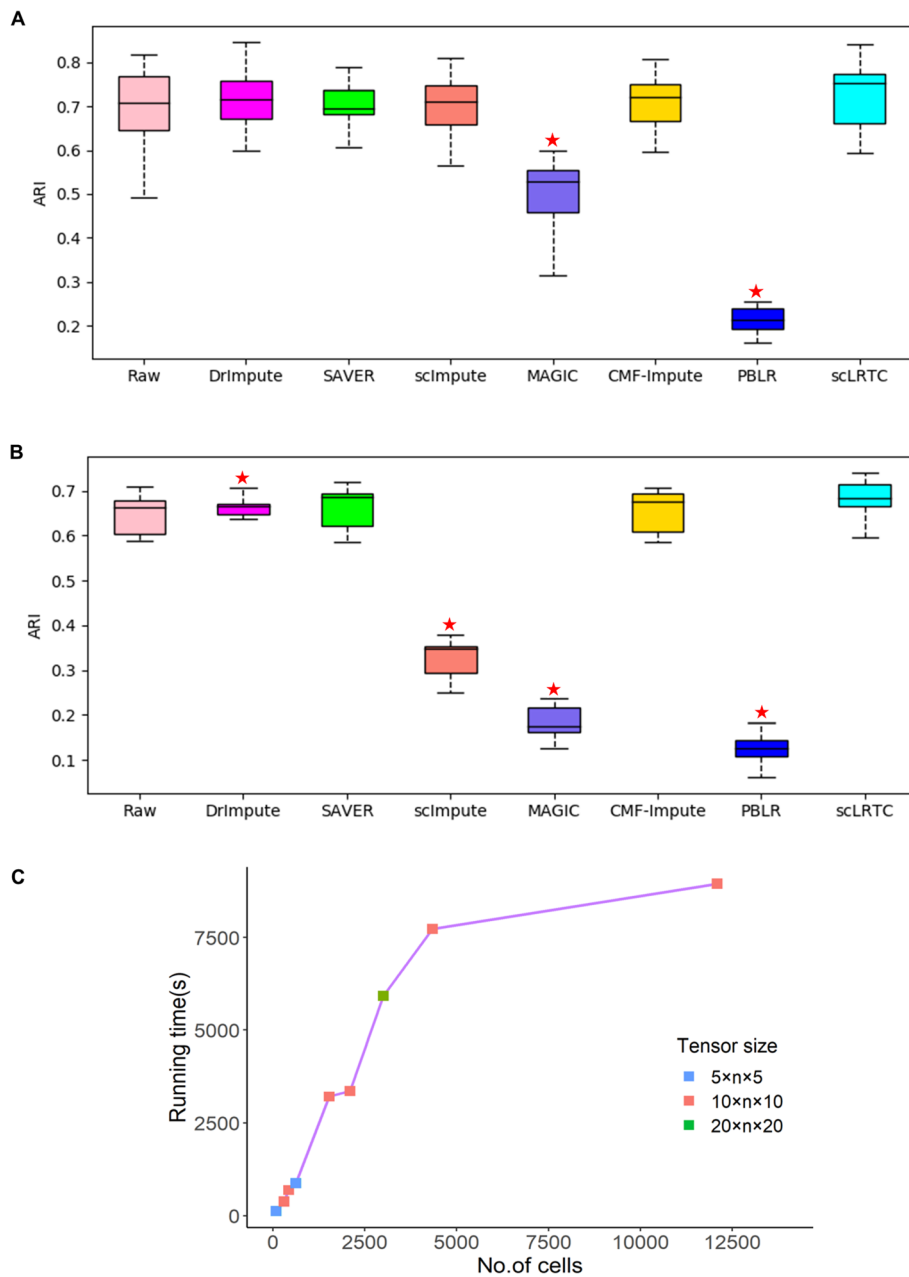
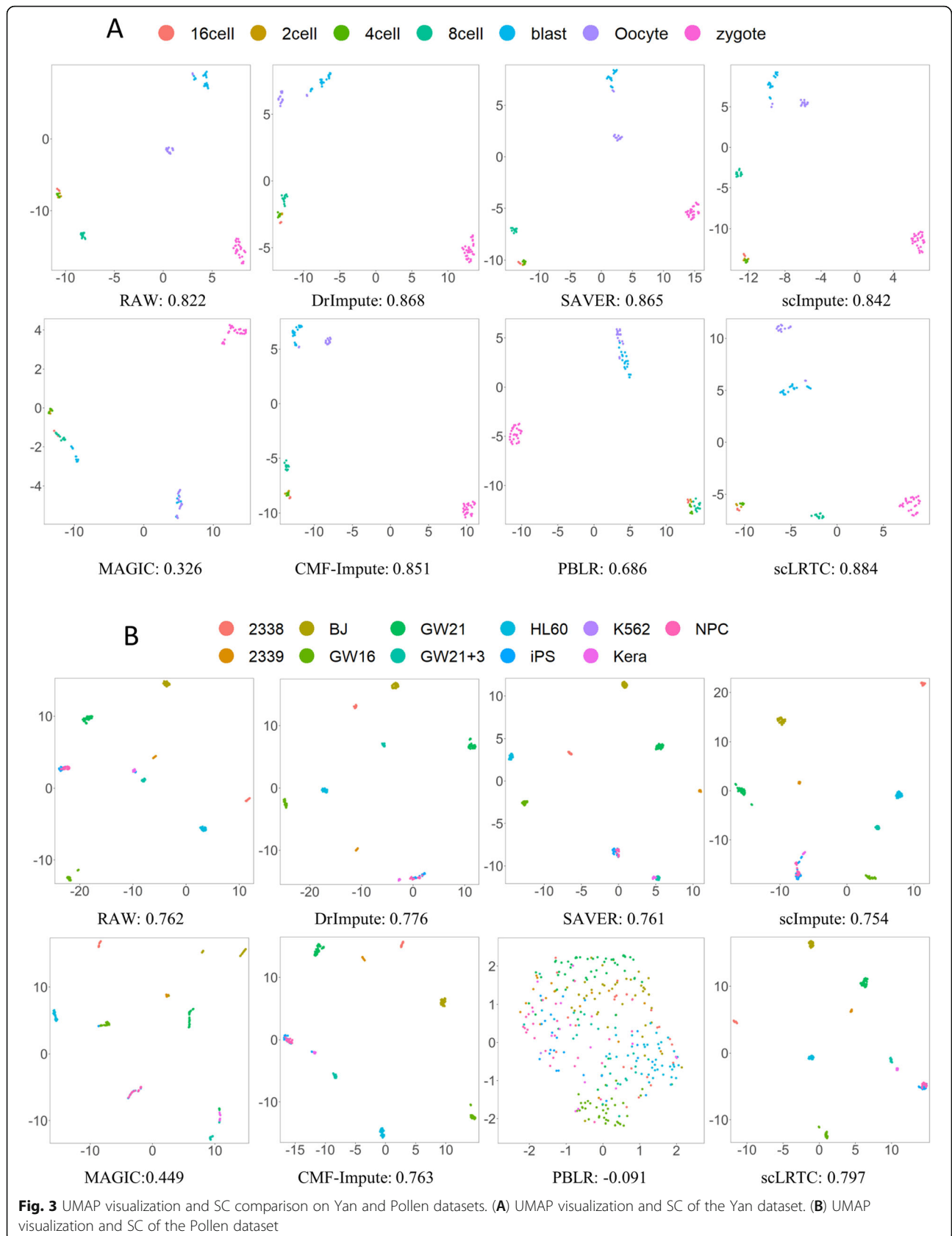


Fig. 2 Performance analysis and comparison of scLRTC and other methods by the t-SNE + K-means clustering. **(A)** ARI obtained by t-SNE + K-means clustering on the Pollen dataset using different algorithms. **(B)** ARI obtained by t-SNE + K-means clustering on the Usoskin dataset using different algorithms. In **(A)** and **(B)**, asterisk indicates the statistically significant difference ($P < 0.05$) between scLRTC and the imputation method of interest using the Wilcoxon rank-sum test. **(C)** Running time of scLRTC for datasets with different sample sizes. The different tensor size setting is represented by different colors, where n is the number of genes

Assessing imputation accuracy through data masking

The data masking evaluation was conducted on the real dataset and simulation datasets. First, 5% of non-zero entries were randomly selected from the Loh dataset, and these values were masked to zeros to generate a new gene expression matrix. Subsequently, seven imputation algorithms were applied for the new gene expression matrix and compared with unmasked data. The

sum of squared errors (SSE) and Pearson correlation coefficient (PCC) between the imputation values and the true values were adopted to evaluate the effect of imputation. Figure 5A presents all the results of the imputation accuracy index of the masked data. The proposed method can recover the missing values with the lowest SSE of 268.8 and the highest PCC of 0.707 in all compared imputation algorithms. Note that the SAVER



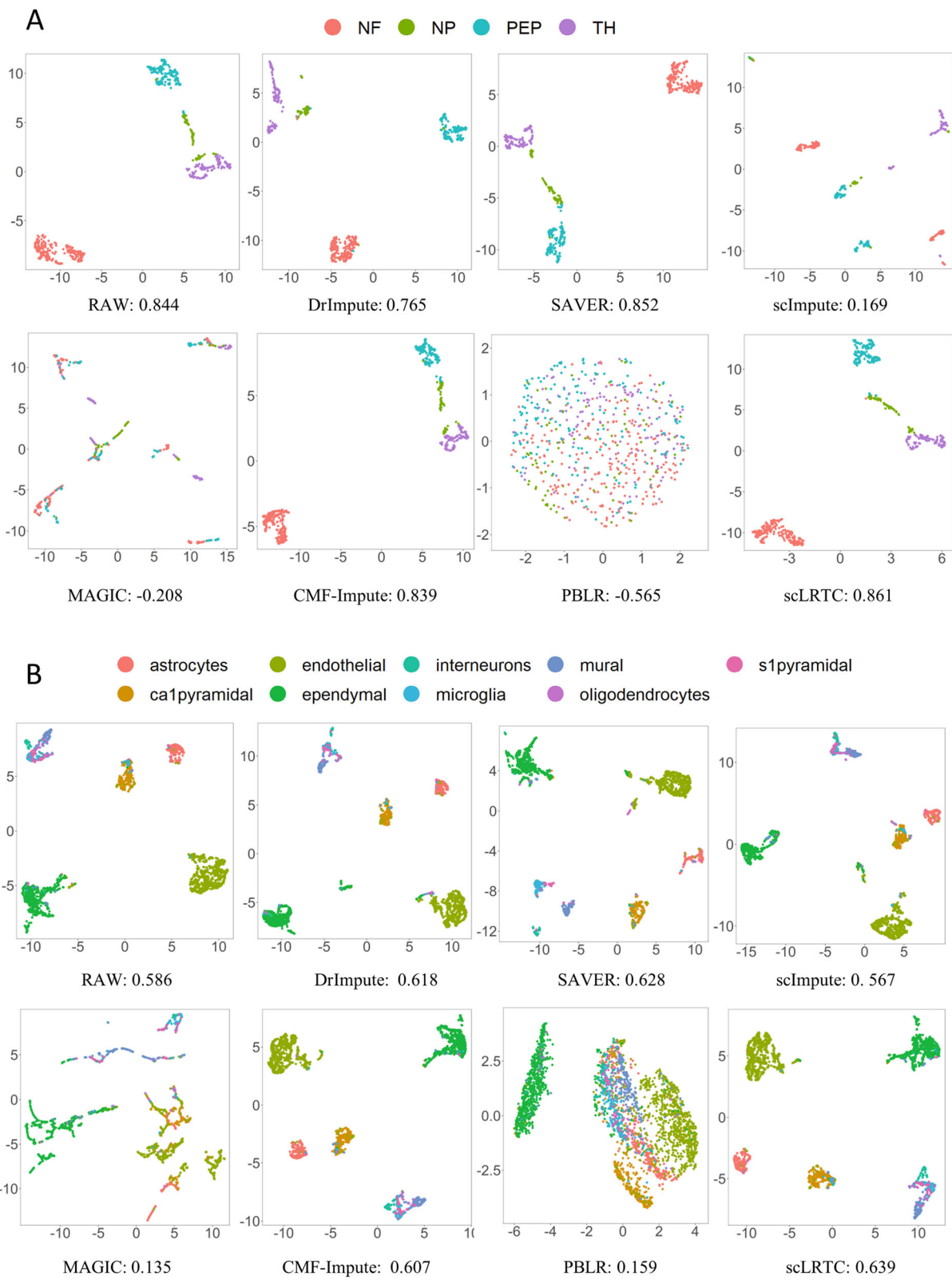
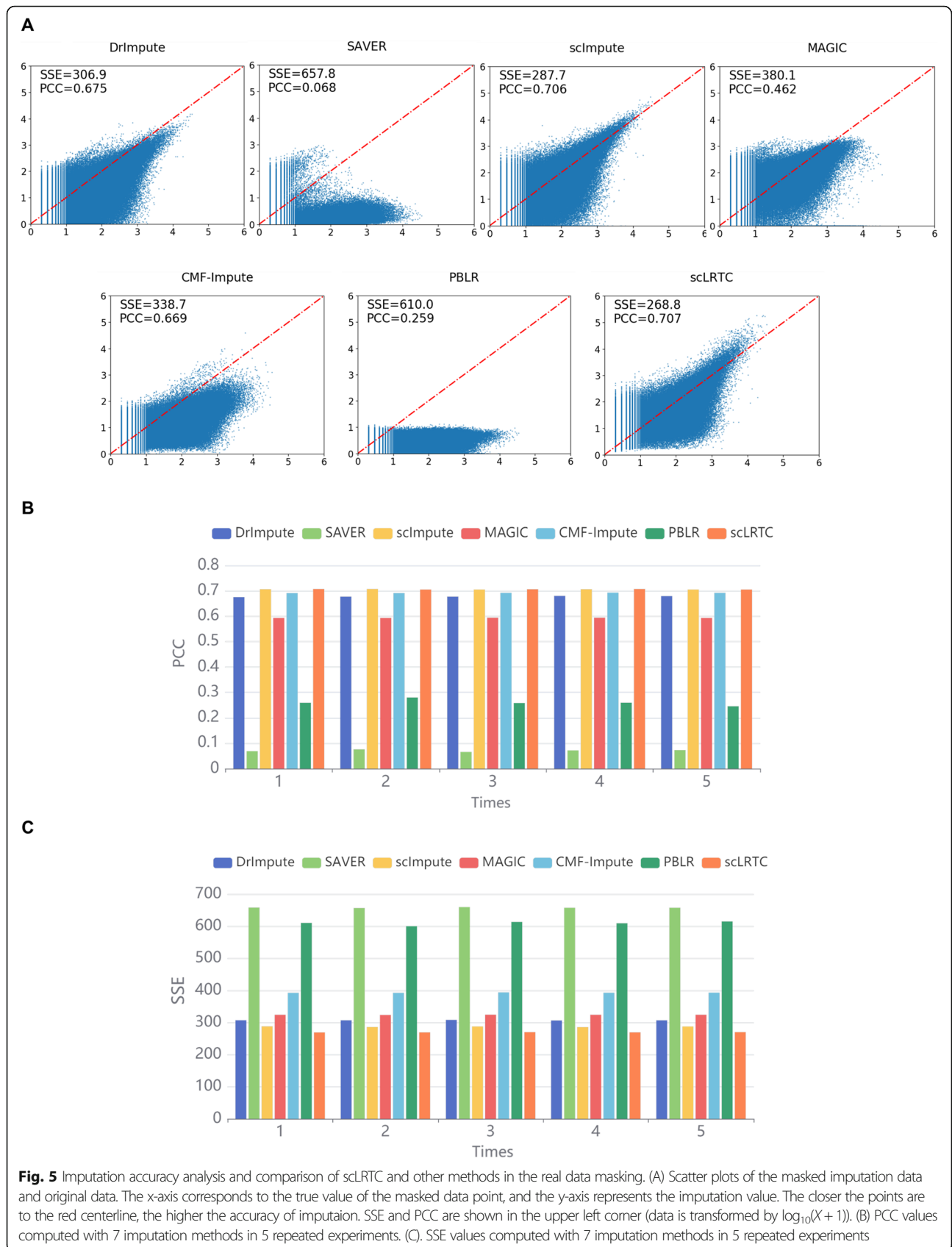


Fig. 4 UMAP visualization and SC comparison on Usoskin and Zeisel datasets. **(A)** UMAP visualization and SC of the Usoskin dataset. **(B)** UMAP visualization and SC of the Zeisel dataset



method persistently underestimates the values, especially among the highly expressed genes. Consistent experimental results were also mentioned in references [3, 15]. To prevent the influence attributed to randomness, we performed 5 masking repetitions for the above experiment. The results of 7 methods in 5 repeated experiments only slightly fluctuate (Fig. 5B and Fig. 5C), demonstrating that the randomness slightly impacts the mentioned results.

Moreover, the performance of the proposed model was tested on single-cell simulation data that involves three cell populations. These data were generated using the Splatter package [29]. Splatter is an R bioconductor package for the reproducible and accurate simulation of scRNA-seq data. We referred to the parameters of simulation dataset provided by CMF-Impute [11] and increased the dropout rate in our experiment. Namely, 40, 50, 60, and 70% of the entries were randomly masked in the expression matrix, corresponding to a shape parameter of dropout logistic function (ds) equaling -0.3, 0, 0.05, and 0.25 respectively. The masked entries were imputed with 7 methods and the imputed results are compared with the real values. Figure 6 shows the visualization results of t-SNE with dropout, unmasked raw data (Full), and 7 imputation methods (including DrImpute, scImpute, MAGIC, SAVER, CMF-Impute, PBLR and scLRRTC) under different dropout rates. It can be seen that the proposed scLRRTC is most consistent with the original data (Full) under the t-SNE visualization, demonstrating that the proposed imputation has a strong ability to restore real cell clusters. Furthermore, we performed the quantitative analysis on the simulation dataset. Figure 7 shows the SSE and PCC values under different dropout rates. With the increase in the dropout rate, the accuracy of all imputation methods is affected. However, the proposed scLRRTC is suggested to exhibit the optimal performance among 7 methods.

Evaluating imputation accuracy through correlation analysis and differential expression analysis

The ability of the imputation method was evaluated to restore gene-gene and cell-cell relationships in complex tissues. The simulated data were employed with a dropout rate of 40% ($ds = -0.3$) to calculate the gene-gene and cell-cell correlation matrix, and $\log_{10}(X + 1)$ was set as the result after imputation. In the cell-to-cell correlation heat map (Fig. 8A), the color of MAGIC and the proposed scLRRTC is the closest to the heat map of Full. For the heat map of gene-gene correlation (Fig. 8B), scImpute and the proposed scLRRTC are the closest ones to the expression heat map of Full in color, while MAGIC deviates the most. And then the violin chart was used to display its expression distribution. We find

the violin chart of scLRRTC is the closest to the unmasked raw data (Full) in the appearance comparison (Fig. 9A), indicating that the position and the upper quartile comply with Full. It is suggested that the data by the scLRRTC complement here achieves the most consistent distribution with that of Full. In summary, the proposed method can effectively restore the true gene-gene and cell-cell relationship.

In addition, it is considered that the imputation method should be capable of recovering true differential expressed (DE) genes and reducing the production of false positive genes. Since gold standard of DE genes has been rarely formulated in real datasets, 6 imputation methods (DrImpute, SAVER, scImpute, MAGIC, CMF-Impute and scLRRTC) were compared for their capabilities to recover DE genes in the simulation data. The differential expression analysis was performed by using the MAST [39], and the true DE genes identified from the complete data were considered the reference. In terms of the respective method, the DE genes were extracted, which are considered significant by controlling P -value < 0.01 and comparing them with the true DE genes. Figure 9B presents the average ROC (Receiver Operating Characteristic) curves of different imputation methods by considering the indices of recall and precision. scLRRTC is found to achieve the highest score (AUC (Area Under the Curve) = 0.971) for detecting DE genes, demonstrating that scLRRTC is valid to recover more DE genes and detect less false-positives genes.

Evaluating imputation accuracy through cell trajectory inference

A common task of single cell RNA sequence analysis is to rebuild the lineage trajectory and infer the differentiation and progenitor status of single cells, which is a research hotspot over the past few years. Besides, a wide range of algorithms have been developed in this field. For instance, TSCAN [40] performed the differential expression and time series analysis on single-cell expression data, which classified individual cells according to the progress of biological processes. However, TSCAN did not perform dropout imputation for the data reprocessing. Thus, in this study, the scLRRTC imputation was integrated into TSCAN, and its performance was compared in the pseudotime inference of the Petropoulos dataset. The Petropoulos data consists of the single cells from five stages of human preimplantation embryonic development from developmental day (E) 3 to day 7. Notably, though the cells at each time point may not be homogeneous, the time label can be exploited to represent an overall developmental trajectory. Accordingly, the known time label acts as the ground truth, and the performance of pseudotime inference is evaluated with TSCAN, as input by the raw data and the imputed data

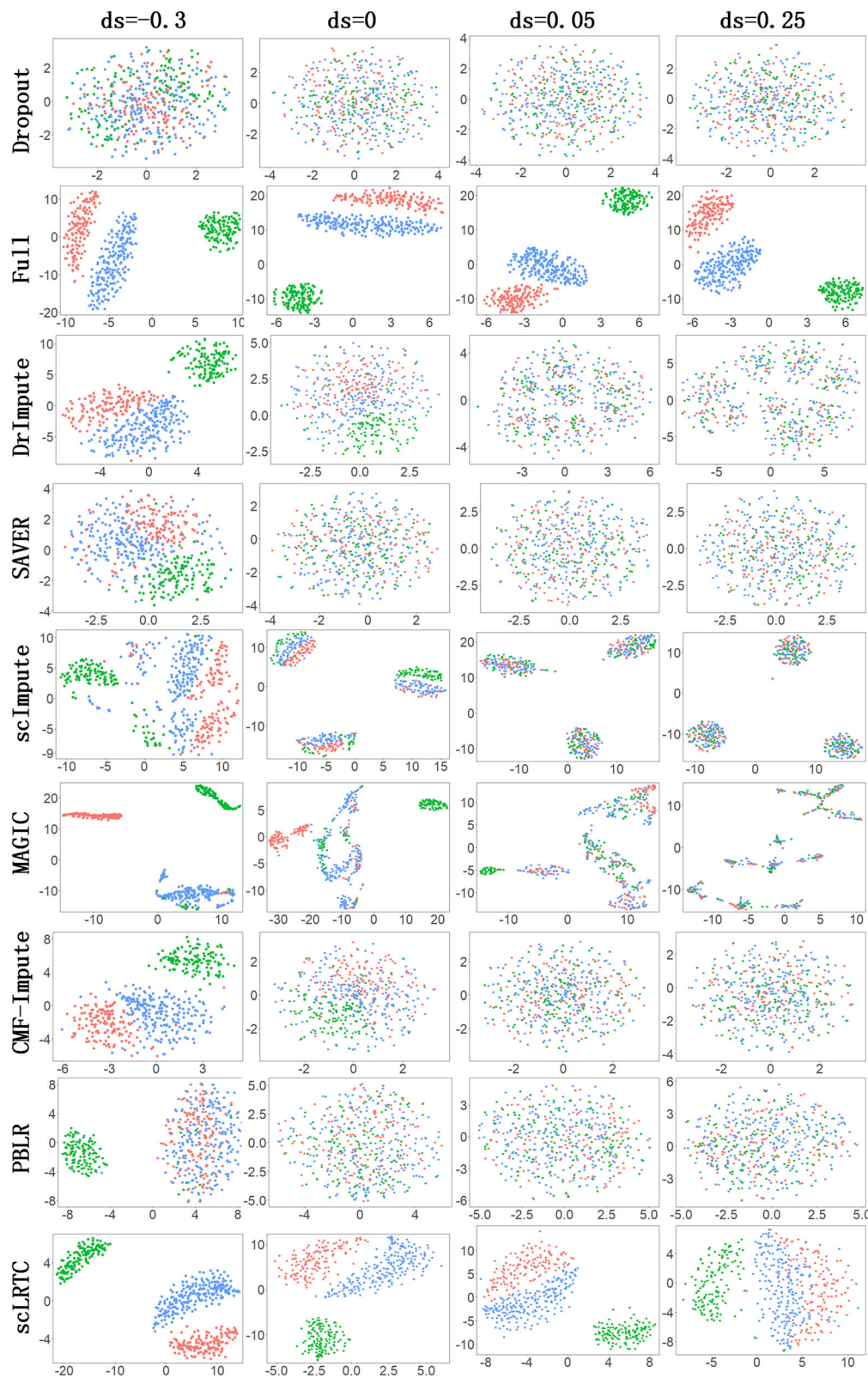
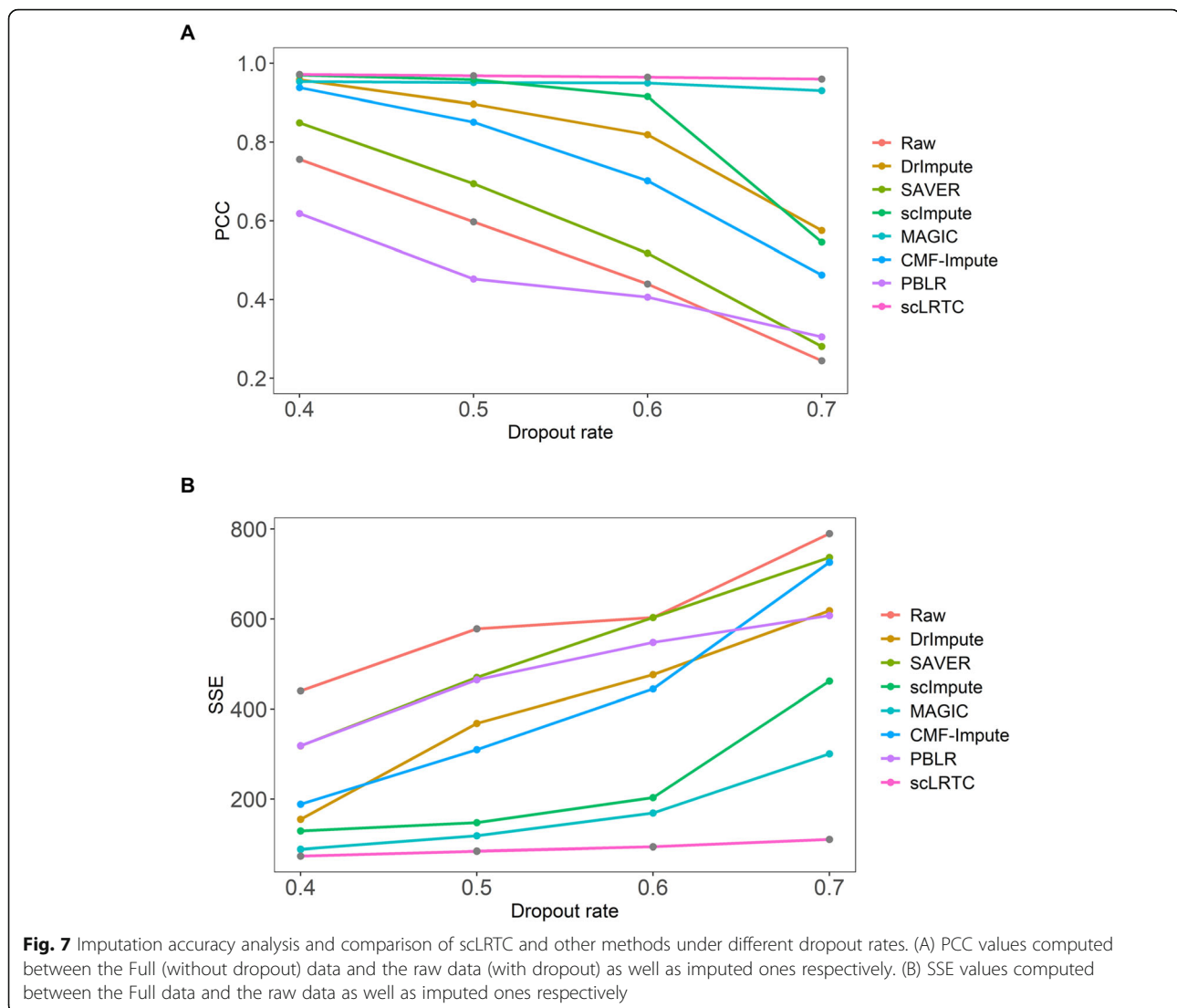


Fig. 6 Visual distinguishability and comparison of scLRTC and other methods on simulation datasets with various dropout rates. We use t-SNE to visualize the cell gene expression matrix, and apply different algorithms for imputation. Each column represents a ds , which controls the ratio of zero



with 6 different methods (scLRTC, SAVER, scImpute, DrImpute, CMF-Impute and MAGIC). Furthermore, Pseudotime ordering score (POS) [40] and Kendall's rank correlation score (KRCS) were used to measure the consistency of time label and pseudotime order derived from the data. The results are presented in Fig. 10.

It is therefore suggested that the proposed method has improved on both POS and KRCS indicators compared with the original data. For the SAVER method, it has an improvement on KRCS, whereas the POS score decreases. In terms of other methods, the pseudotime trajectory by DrImpute and scImpute starts at E3 and ends at E6, the pseudotime trajectory of CMF-Impute starts at E3 and ends at E5, and the trajectory reconstruction error is relatively large. The accuracy of MAGIC from E3 to E5 is relatively high, whereas at E6 and E7 stages, a big discrepancy is identified with the real label, which introduces errors.

Discussion

Since single-cell RNA has a limited extraction efficiency, the occurrence of dropout events adversely affects the downstream analysis. However, the single-cell data imputation is not explicitly involved in the most used scRNA-seq tools for cell clustering, dimensionality reduction visualization, cell type recognition and lineage reconstruction, so it is of high research significance. In this paper, a novel low-rank tensor completion method (scLRTC) is proposed to impute the scRNA sequence data where dropout is present. ScLRTC, a data-driven method, fully considers the similarity and heterogeneity between cells. It builds a third-order tensor representation and employs a low-rank tensor completion model by adopting the ADMM algorithm to achieve the data imputation. This study also inputs the data with scLRTC imputation into SC3 clustering and carries out the clustering by first conducting t-SNE dimensionality

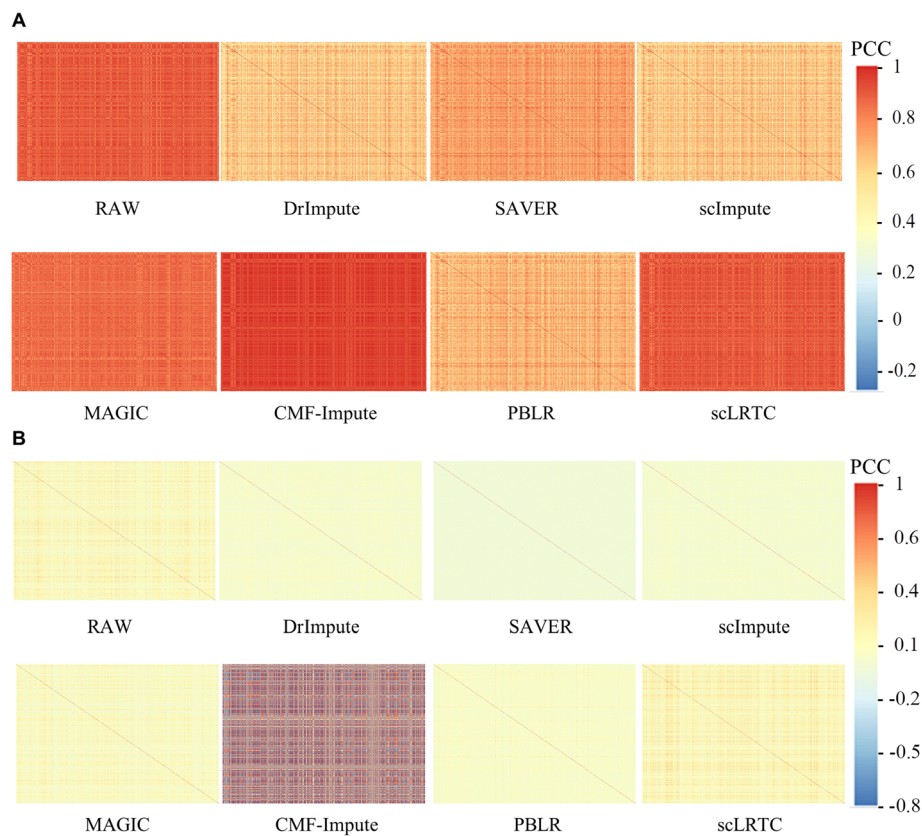


Fig. 8 Correlation analysis and comparison of scLRTC and other methods. The more similar the heat map is to the raw heat map, the better the imputation effect. **(A)** Visualized heat map of cell-cell correlation matrix. **(B)** Visualized heat map of gene-gene correlation matrix

reduction and then implementing K-means. Accordingly, it is reported that scLRTC is capable of increasing the clustering accuracy of real data at different dropout rates, as well as improving the quality of cell visualization. Moreover, by integrating the proposed scLRTC into TSCAN, we find it improves the accuracy of pedigree reconstruction and pseudotime inference.

Although we have verified that our scLRTC is superior to other popular methods on some real and simulation datasets, we cannot guarantee it is superior to all other methods on all datasets. We found our method is effective in imputing the scRNA-seq dataset with a high missing rate, mainly based on the following aspects. (1). Compared to ScImpute, PBLR and other clustering based methods, when the data missing rate is high, incorrect clustering result will affect the subsequent imputation process. Our scLRTC makes full use of the cell similarity to construct a low-rank tensor, which can reduce the impact of highly missing data on the imputation process. (2). For SAVER and other methods based on the statistical model, they normally impute the entire data under a given data distribution assumption. When the data distribution does not meet this assumption, the completion effect will be affected. But imputation by the low rank tensor of scRNA-seq data can avoid the influence

of data distribution assumption. (3). For the scRNA-seq data, although the data itself has redundancy, the rank estimation of the original gene expression matrix is easy to be affected when the data has a high missing rate. Whereas, the rank estimation of the tensor constructed in our scLRTC can be tracked by the tensor trace norm, which can guarantee the final completion result.

In general, the proposed imputation method can be regarded as one powerful complement to current scRNA sequence data analysis. Our tensor based imputation algorithm can be further improved in the future work. For instance, because of the tensor model in scLRTC is relatively independent, we will develop the single-cell completion based on the parallel computing to improve the time complexity of scLRTC. Besides, we currently only use the similarity between cells to build a low-rank tensor. We can also consider the similarity between genes, and combine the similarity between cells and genes to build a higher-order tensor, and then complete the imputation under the tolerable computational complexity. In addition, we currently developed the scLRTC based on MATLAB mainly because there are the tensor related packages so that we can quickly develop our algorithm and verify it in the experiment analysis. In our future

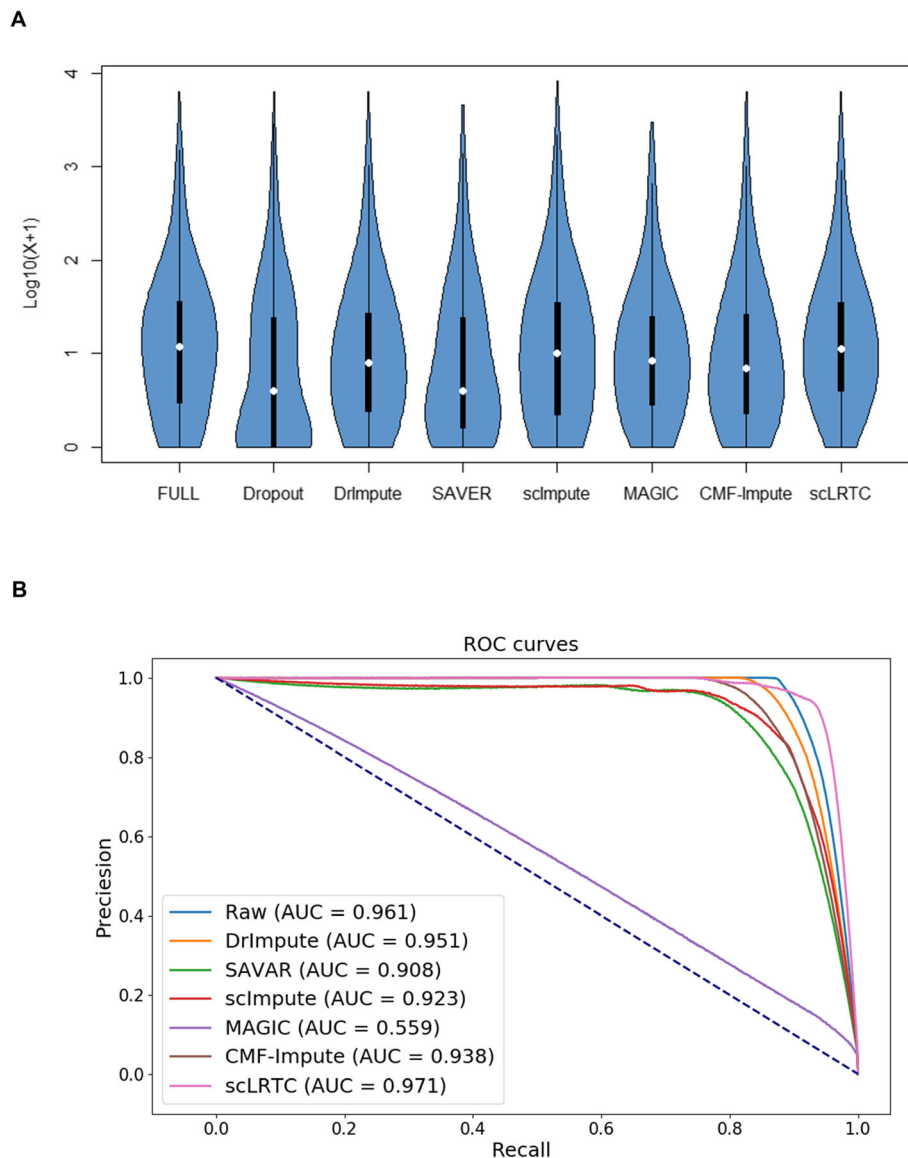


Fig. 9 Violin chart of data expression distribution and accuracy measurements of DE genes by scLRTC and other methods. **(A)** Violin chart of data expression distribution after imputation when the dropout rate is 40% (data is transformed by $\log_{10}(X+1)$). The more similar the shape of the violin is to FULL, the more effective the imputation effect. **(B)** ROC curves and AUC scores of DE genes with different imputation methods. AUC combines the recall rate and precision rate, and the value closer to 1 indicates a better imputation method. Here, the recall rate is defined as the number of true positives divided by the total number of samples that actually belong to the positive class, and the precision rate is the number of true positives divided by the total number of samples labelled as belonging to the positive class

work, we will use R or Python to realize the scLRTC algorithm for providing the widespread use in the bioinformatics community.

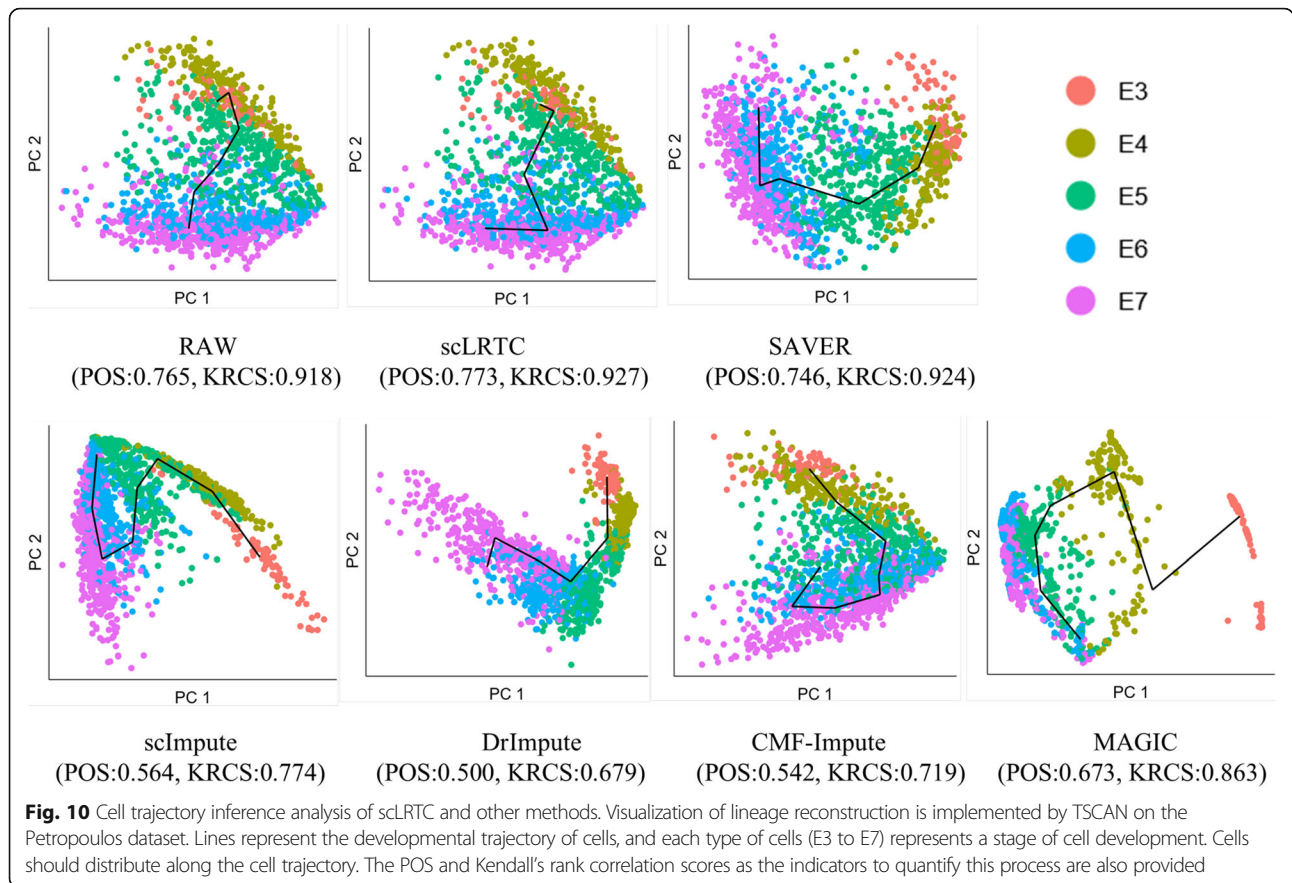
Conclusions

Imputation is an essential step in the use of scRNA-seq. In this work we introduced a novel low-rank tensor completion-based method, termed as scLRTC. Experiments on simulation data and real data sets showed scLRTC to be highly accurate in imputation.

Methods

Datasets

Nine scRNA-seq datasets (i.e., Pollen, Usoskin, Yan, Zeisel, Mouse, PBMC, Chen, Loh and Petropoulos) with different data sizes are used to test the validity of the proposed scLRTC in imputing dropout events. Besides, these datasets fall to three levels (i.e., gold, silver and copper) based on the supporting evidence of cell markers. To be specific, Pollen, Loh, Yan, Zeisel and Mouse datasets are defined as gold standard datasets, in



which all cell markers are defined by complying with experimental conditions or cell lines. Usoskin, PBMC and Chen datasets are defined as the silver standard, with cell markers are calculated and assigned by drawing upon the authors' knowledge of the underlying biology. Petropoulos is considered the copper standard since the cells involved are in the developmental stage (time labeled). Though single cell populations from different time points usually exhibit different expression patterns and biological characteristics, it remains

infeasible to separate different populations at each time point based on time tags alone. Table 3 briefs these scRNA sequence datasets with sizes ranging from 90 (Yan) to 12,089 (Chen), and the number of cell clusters ranges from 4 (Usoskin) to 46 (Chen). Note that the first seven datasets are normally used for the cluster analysis. Furthermore, the first three datasets are from the low-throughput data sequencing platform, and the last four datasets originate from the high-throughput data sequencing platform. Loh is

Table 3 A summary of nine real scRNA-seq datasets used in our experiment

Dataset	Number of clusters	Number of cells	Number of genes	Standard
Pollen	11	301	23,730	Gold
Usoskin	4	622	25,334	Silver
Yan	7	90	20,214	Gold
Zeisel	9	3005	19,972	Gold
Mouse	16	2100	20,670	Gold
PBMC	8	4340	33,694	Silver
Chen	46	12,089	23,284	Silver
Loh	8	429	23,794	Gold
Petropoulos	5	1529	21,749	Copper

employed for the data masking evaluation, and Petropoulos is used for the trajectory reconstruction analysis.

Data preprocessing and normalization

In terms of a given scRNA-seq dataset, its gene expression matrix is recorded as X^c . To reduce the effect of underexpressed genes, the gene expressed in less than or equal to 3 cells is removed [41]. To express the filtered matrix by X^N , a matrix X is then made by taking the \log_2 transformation with a pseudo count 1

$$\begin{aligned} X_{ij} &= \log_2(X_{ij}^N + 1); i = 1, 2, \dots, M, j \\ &= 1, 2, \dots, N \end{aligned} \quad (1)$$

where M denotes the overall number of genes; N is the total number of cells. The pseudo-count is added to avoid infinite values in the parameter estimation in the subsequent data analysis. The logarithmic transformation has an advantage that it can prevent a small number of large observation values from being significantly affected in the data imputation.

Tensor based model for scRNA-seq data imputation

Single-cell dropout events can be formulated as a missing value estimation problem. The core problem of missing value estimation refers to how to develop the relationship between known elements and unknown elements. The scRNA-seq data usually consist of single cells from different cell types, and single cells exhibiting the identical type have similar expression patterns. For this reason, it is assumed that the basic true expression of scRNA-seq data can be approximatively considered as a low-rank matrix. The low-rank matrix restoration essentially complies with the correlation between the rows and columns of a matrix, therefore creates a direct and effective imputation strategy.

A recent study suggested that taking full advantage of the presence of low-rank submatrix can improve the imputation performance compared to traditional low-rank matrix recovery methods [14, 15]. However, the low-rank submatrix constructed by clustering is easily influenced by the clustering effect, and the low-rank tensor can be constructed to capture more correlations of similar single cell compared to the low-rank submatrix form. Based on this motivation, the two-dimensional low-rank matrix is extended to the three-order low-rank tensor with the high correlation of scRNA-seq data. Besides, a novel low-rank tensor model is built for single-cell gene expression data, and the tensor trace norm [17] is employed to approximate the rank of the tensor, finally the missing data are rebuilt and the cell's gene expression is restored.

Tensor construction

We construct the tensor form of single-cell which fully considers the high correlation of scRNA-seq data. Specific to a given cell X_i , the Pearson correlation coefficient between cells is first calculated and sorted in a descending order. Subsequently, the gene expression of cell X_i and its $K-1$ cells with the highest correlation are adopted to build a matrix $Mat_i \in \mathbb{R}^{K \times M}$, where M denotes the number of genes. Subsequently, the difference between cell X_i and other cells are measured in the matrix Mat_i . To be specific, the Euclidean distance is calculated from all cells in Mat_i to X_i and sorted in an ascending order as $D_i = (d_{1i}, d_{2i}, \dots, d_{Ki})$. Next, the direction (angle) similarity of cell gene expression is measured by calculating the cosine similarity from all cells in Mat_i to X_i , and then it is recorded in an ascending order as $C_i = (c_{1i}, c_{2i}, \dots, c_{Ki})$. Lastly, the similarity between cells is measured according to the absolute value of the difference of cell gene expressions, i.e., the Chebyshev distance from all cells in Mat_i to X_i is calculated and restored in an ascending order as $Q_i = (q_{1i}, q_{2i}, \dots, q_{Ki})$.

Three distance vectors obtained from X_i are combined into a feature vector $Vec_i = \{D_i, C_i, Q_i\}$ with a size of $3K \times 1$. Likewise, the feature vector Vec_j can be obtained from the other cell X_j . By calculating the distance between two feature vectors, the $P-1$ Vec_j closest to Vec_i is searched, and these matrices are merged to build a third-order tensor $\mathcal{Y} \in \mathbb{R}^{K \times M \times P}$ for the cell X_i (as shown in Fig. 11).

Our method fully exploits the similarity between cells to construct a low-rank tensor, which can avoid the imputation influence by the high missing rate of scRNA-seq data in the completion process. On the other hand, the tensor trace norm is used to track the rank of tensor and solved by the ADMM algorithm, which can guarantee the imputation results more accurate and robust compared with other heuristic tensor completion methods (Tucker, Parafac and SVD).

Tensor fold and unfold

During the tensor analysis, it is convenient to unfold a tensor into a matrix. The "unfold" operation along the k th mode on a tensor \mathcal{Y} is defined as

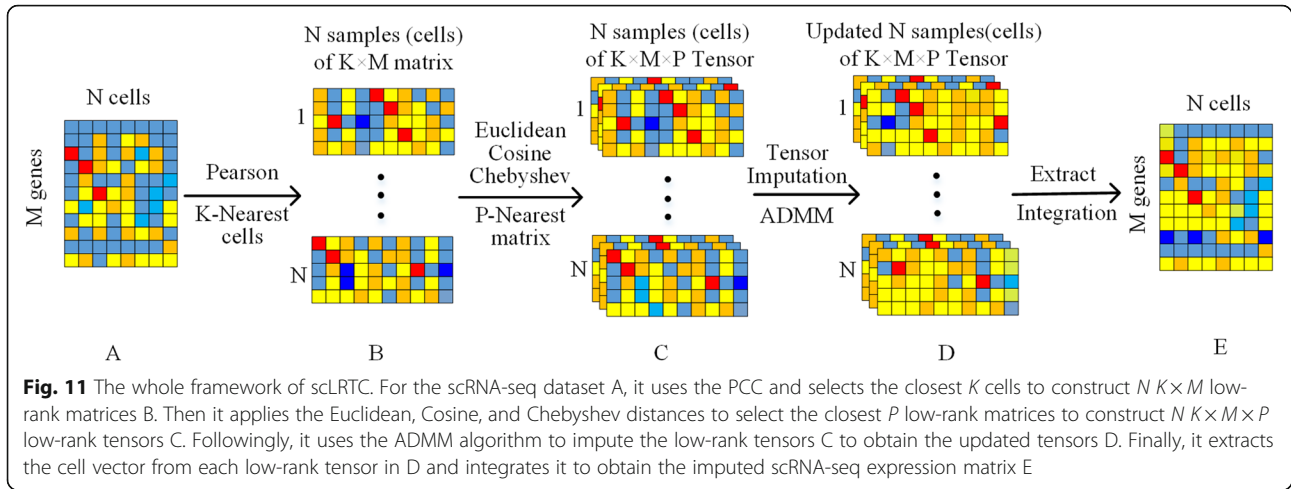
$$\text{unfold}_k(\mathcal{Y}) = \mathcal{Y}_{(k)} \in \mathbb{R}^{I_k \times (I_1 \dots I_{k-1} I_{k+1} \dots I_n)} \quad (2)$$

The opposite operation "fold" is defined as

$$\text{fold}_k(\mathcal{Y}_{(k)}) = \mathcal{Y} \quad (3)$$

Tucker decomposition and denoising

Tensor Tucker decomposition is recognized as a form of high-order principal component analysis. The HOSVD



method [42] is available for decomposing a third-order tensor $\mathcal{Y} \in \mathbb{R}^{I \times J \times K}$ by

$$\mathcal{Y} \approx \mathcal{G} \times_1 A \times_2 B \times_3 C \tag{4}$$

where $A \in \mathbb{R}^{I \times P}$, $B \in \mathbb{R}^{J \times Q}$ and $C \in \mathbb{R}^{K \times R}$ denote factor matrices, which can be considered the main components of the corresponding mode. The tensor $\mathcal{G} \in \mathbb{R}^{P \times Q \times R}$ refers to the core tensor, representing the level of interaction between different components.

Next, a hard threshold function is set for factor matrices A , B , C to eliminate the effect of some low value components after the Tucker decomposition. Subsequently, the third-order tensor $\hat{\mathcal{Y}}$ is restored by updating $\hat{\mathcal{G}}$, i.e., the convergence of the current tensor data is ensured by the iterative computation, and the denoising effect of some mutation elements is achieved in the tensor form of scRNA-seq data.

Tensor trace norm

The trace norm of a tensor is defined as [17].

$$\|\mathcal{Y}\|_* = \sum_{i=1}^n \alpha_i \|\mathcal{Y}_{(i)}\|_* \tag{5}$$

where $\alpha_i \geq 0$, $\sum_{i=1}^n \alpha_i = 1$, $\mathcal{Y}_{(i)}$ expresses the matrix along the i th mode. In fact, the trace norm of a tensor refers to a convex combination of the trace norms of all matrices expanded along the respective mode. Notably, when n is equal to 2 (the matrix case), the definition of the tensor's trace norm complies with the matrix case.

Tensor completion model

By using M_i to replace \mathcal{Y} , the tensor completion model is expressed as

$$\begin{aligned} \min_{\mathcal{Y}, M_1, \dots, M_n} \sum_{i=1}^n \alpha_i \|M_{(i)}\|_* \\ \text{s.t. } \mathcal{Y}_{\Omega} = \mathcal{T}_{\Omega} \\ \mathcal{Y} = M_i, i = 1, \dots, n. \end{aligned} \tag{6}$$

where α_i denotes the coefficient; $M_{(i)}$ represents the unfold matrix of the tensor along the i th mode; \mathcal{T} is the known tensor; \mathcal{Y} expresses the reconstructed tensor; Ω in \mathcal{T}_{Ω} is the index of non-zero observation value.

Tensor model solution

The mentioned model can be solved by adopting the alternating direction method of multipliers (ADMM) algorithm. The augmented Lagrangian function is defined as

$$\begin{aligned} L_{\rho}(\mathcal{Y}, M_1, \dots, M_n, y_1, \dots, y_n) \\ = \sum_{i=1}^n \alpha_i \|M_{(i)}\|_* + \langle \mathcal{Y} - M_i, y_i \rangle + \\ \frac{\rho}{2} \|\mathcal{Y} - M_i\|_F^2 \end{aligned} \tag{7}$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product; $\|\cdot\|_F^2$ represents the F-norm, i.e., the root of the square sum of all elements; y_i is the Lagrange multiplier; ρ expresses the penalty parameter.

According to the framework of ADMM, M_i, \mathcal{Y}, y_i can be iteratively updated as

$$\{M_1^{k+1}, \dots, M_n^{k+1}\} = \underset{M_1, \dots, M_n}{\operatorname{argmin}} L_{\rho}(\mathcal{Y}^k, M_1, \dots, M_n, y_1^k, \dots, y_n^k) \tag{I}$$

$$\mathcal{Y}^{k+1} = \underset{\mathcal{Y}}{\operatorname{argmin}} L_{\rho}(\mathcal{Y}, M_1^{k+1}, \dots, M_n^{k+1}, y_1^k, \dots, y_n^k) \tag{II}$$

$$y_i^{k+1} = y_i^k - \rho(M_i^{k+1} - \mathcal{Y}^{k+1}) \tag{III}$$

From the augmented Lagrangian function in 1), $\mathcal{Y}_{(i)}^k$, y_i^k is fixed and minimized to yield

$$\begin{aligned}
M_{i(i)}^{k+1} &= \arg \min_{M_{i(i)}} \alpha_i \|M_{i(i)}\|_* \\
&\quad + \frac{\rho}{2} \langle M_{i(i)}, M_{i(i)} \rangle - \rho \left\langle M_{i(i)}, \mathcal{Y}_{(i)}^k + \frac{1}{\rho} y_{i(i)}^k \right\rangle \\
&= \arg \min_{M_{i(i)}} \alpha_i \|M_{i(i)}\|_* + \frac{\rho}{2} \left\| M_{i(i)} - \mathcal{Y}_{(i)}^k - \frac{1}{\rho} y_{i(i)}^k \right\|_F^2 \\
&\quad - \left\langle \mathcal{Y}_{(i)}^k + \frac{1}{\rho} y_{i(i)}^k, \mathcal{Y}_{(i)}^k + \frac{1}{\rho} y_{i(i)}^k \right\rangle \quad (8)
\end{aligned}$$

Thus, the optimal solution of $M_{i(i)}^{k+1}$ is

$$\begin{aligned}
M_{i(i)}^{k+1} &= \arg \min_{M_{i(i)}} \alpha_i \|M_{i(i)}\|_* \\
&\quad + \frac{\rho}{2} \left\| M_{i(i)} - \left(\mathcal{Y}_{(i)}^k + \frac{1}{\rho} y_{i(i)}^k \right) \right\|_F^2 \quad (9)
\end{aligned}$$

The above Eq. (9) is proven to generate a closed-form in recent references [43, 44], so it can be solved by calculating the singular value thresholding operator $D_\tau(\cdot)$. In terms of any matrix X , the singular value decomposition (SVD) is performed to obtain $X = U\Sigma V^T$, where U , V are orthogonal singular vectors, and $\Sigma \in \mathbb{R}^{r \times r}$ comprises singular values $\sigma_1, \dots, \sigma_r$, $r = \min\{m, n\}$. The singular value thresholding operator can be defined as $D_\tau(X) = U\Sigma_\tau V^T$, where $\Sigma_\tau = \text{diag}(\max(\sigma_i - \tau, 0))$. Thus, it yields

$$M_{i(i)} = D_{\frac{\alpha_i}{\rho}} \left(\mathcal{Y}_{(i)} + \frac{1}{\rho} y_{i(i)} \right) \quad (10)$$

By folding $M_{i(i)}$ to get

$$M_i = \text{fold}_i(M_{i(i)}) \quad (11)$$

From the augmented Lagrangian function in 2), it can be minimized by fixing M_i^k and y_i^k , and the optimal solution is obtained as

$$\begin{aligned}
\mathcal{Y}^{k+1} &= \arg \min_{\mathcal{Y}} \sum_{i=1}^n \langle \mathcal{Y}, y_i^k \rangle + \frac{\rho}{2} \langle \mathcal{Y} - M_i^{k+1}, \mathcal{Y} - M_i^{k+1} \rangle \\
&= \arg \min_{\mathcal{Y}} \sum_{i=1}^n \frac{\rho}{2} \langle \mathcal{Y}, \mathcal{Y} \rangle + \langle \mathcal{Y}, \rho M_i^{k+1} - y_i^k \rangle \quad (12)
\end{aligned}$$

Take the derivative of (12) with respect to \mathcal{Y} and set it equal to 0 to yield

$$n\rho\mathcal{Y} - \left(\sum_{i=1}^n \rho M_i^{k+1} - y_i^k \right) = 0 \quad (13)$$

So

$$\mathcal{Y}^{k+1} = \frac{1}{n} \left(\sum_{i=1}^n M_i^{k+1} - \frac{1}{\rho} y_i^k \right) \quad (14)$$

For $\mathcal{Y} = \mathcal{Y}_{\bar{\Omega}} + \mathcal{Y}_{\Omega}$, $\mathcal{Y}_{\Omega} = \mathcal{T}_{\Omega}$ is known, so only $\mathcal{Y}_{\bar{\Omega}}$ is updated

$$\mathcal{Y}_{\bar{\Omega}}^{k+1} = \frac{1}{n} \left(\sum_{i=1}^n M_i^{k+1} - \frac{1}{\rho} y_i^k \right)_{\bar{\Omega}} \quad (15)$$

After the tensor completion reconstruction, the corresponding cell expression $X_i (i = 1, 2, \dots, N)$ is selected from the respective tensor model, and the gene information of each cell is restored, and a matrix X^* representing the complete scRNA-seq data is lastly formed. Since the gene expression of the cell is non-negative, the matrix X^P is defined after the imputation

$$X^P = \begin{cases} X_{ij}^P = 0 & \text{if } (X_{ij}^* < 0) \\ X_{ij}^P = X_{ij}^* & \text{others} \end{cases} \quad (16)$$

In brief, the entire scLRTC algorithm process is expressed as.

Input: X, ρ, a and K

Output: \hat{X} (X after imputing)

- 1: **for** $m = 1$ to N
 - 2: Construction of low rank tensor T_{Ω} of X_m
 - 3: Set $\mathcal{Y}_{\Omega} = T_{\Omega}$ and $\mathcal{Y}_{\bar{\Omega}} = 0, \mathcal{Y} = M_i$
 - 4: **for** $k = 0$ to K **do**
 - 5: **for** $i = 1$ to n **do**
 - 6: $M_{i(i)} = D_{\frac{\alpha_i}{\rho}}(\mathcal{Y}_{(i)} + \frac{1}{\rho} y_{i(i)})$
 - 7: $M_i = \text{fold}_i[M_{i(i)}]$
 - 8: **end for**
 - 9: $\mathcal{Y}_{\bar{\Omega}} = \frac{1}{n} (\sum_{i=1}^n M_i - \frac{1}{\rho} y_i)_{\bar{\Omega}}$
 - 10: $y_i = y_i - \rho(M_i - \mathcal{Y})$
 - 11: **end for**
 - 12: extract \hat{X}_m from \mathcal{Y}
 - 13: **end for**
 - 14: $X^* = (\hat{X}_1, \hat{X}_2, \dots, \hat{X}_m)$ and use Eq. (16) to remove negative values
-

Evaluation measures

To objectively evaluate the effectiveness of the proposed low-rank tensor completion method for single-cell RNA-seq data, the reconstructed data are used for the cell clustering, and two clustering indicators with the normalized mutual information (NMI) and the adjusted rand index (ARI) are adopted to quantify the consistency between inferred and predefined cell clusters in the respective scRNA-seq data. Subsequently, the silhouette coefficient (SC) is adopted to assess the visual effect of cell dimensionality reduction. Lastly, Pseudotemporal ordering score (POS) and KRCS are used to evaluate the accuracy of cell trajectory analysis and imputation.

Denote that $U = \{\mu_1, \mu_2, \dots, \mu_P\}$ is adopted to represent the true partition of P classes, $V = \{v_1, v_2, \dots, v_K\}$ is used to denote the partition given by clustering results, n_i and n_j are represented as the number of the class μ_i and cluster v_j , respectively, and n_{ij} is expressed as the number of observations in both class μ_i and cluster v_j .

ARI is then formulated as:

$$\frac{\sum_{i=1}^P \sum_{j=1}^K \binom{n_{ij}}{2} - \left[\sum_{i=1}^P \binom{n_i}{2} \sum_{j=1}^K \binom{n_j}{2} \right] / \binom{n}{2}}{\frac{1}{2} \left[\sum_{i=1}^P \binom{n_i}{2} \sum_{j=1}^K \binom{n_j}{2} \right] - \left[\sum_{i=1}^P \binom{n_i}{2} \sum_{j=1}^K \binom{n_j}{2} \right] / \binom{n}{2}} \quad (17)$$

where $n = \sum_{i=1}^P n_i = \sum_{j=1}^K n_j$.
NMI is expressed as

$$NMI = \frac{2I(U, V)}{H(U) + H(V)} \quad (18)$$

where $I(U, V)$ expresses the amount of mutual information between U and V

$$I(U, V) = \sum_{i=1}^P \sum_{j=1}^K \frac{|u_i \cap v_j|}{N} \log \frac{N |u_i \cap v_j|}{|u_i| \times |v_j|} \quad (19)$$

$H(U)$ and $H(V)$ are the entropies of partitions U and V

$$\begin{aligned} H(U) &= - \sum_{i=1}^P \frac{u_i}{N} \log \frac{u_i}{N}, H(V) \\ &= - \sum_{i=1}^K \frac{v_i}{N} \log \frac{v_i}{N} \end{aligned} \quad (20)$$

where N is the total number of cells.

SSE is written as

$$\text{sqrt} \left(\sum_{i=1}^n \sum_{j=1}^n (X_{ij} - P_{ij})^2 \right) \quad (21)$$

where X_{ij} denotes the true gene expression; P_{ij} represents the predicted gene expression.

SC is expressed as

$$SC = \text{average} \left(\sum_{i=1}^N \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \right) \quad (22)$$

where i denotes the i th cell, $a(i) = \text{average}$ (i to all other cells in the cluster to which it belongs), $b(i) = \text{min}$ (the average distance from i to all cells in the other cluster).

Abbreviations

scRNA-seq: Single-cell RNA-sequencing; PCC: Pearson correlation coefficient; SSE: sum of squares of error; ARI: Adjusted Rand Index; NMI: Normalized Mutual Information; SC: silhouette coefficient; POS: Pseudotemporal ordering score; KRCS: Kendall's rank correlation score; ADMM: alternating direction method of multipliers

Authors' contributions

Z.L. provided the methodology. X.P., Z.L. and S.Q. designed the algorithm. X.P., M.Y. and H.H. arranged the datasets and performed the analysis. X.P. drafted the manuscript. Z.L. and S.Q. reviewed and edited the manuscript. All authors read and approved the final manuscript.

Funding

This work was supported by the National Natural Science Foundation of China under Grant No. 12171434, Zhejiang Provincial Natural Science Foundation of China under Grant No. LZ19A010002.

Availability of data and materials

There are no new data associated with this article. Published datasets used in this study: Pollen et al.'s dataset ([20], P041736) of human tissues cells is available at <https://www.ncbi.nlm.nih.gov/sra?term=SRP041736>. Usoskin et al.'s dataset ([21], GSE59739) of mouse lumbar dorsal root ganglion cells is available at <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE59739>. Yan et al.'s dataset ([22], GSE36552) of the human embryos cells is available at <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE36552>. Zeisel et al.'s dataset ([23], GSE60361) of the mouse cortex and hippocampus cells is available at <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE60361>. The mouse bladder cell dataset [24] originates from the mouse cell atlas project is available at <https://figshare.com/s/865e694ad06d5857db4b>. From the raw count matrix, 2100 cells are selected from the bladder tissue. PBMC dataset [25] is downloaded from the 10X genomics website (<https://support.10xgenomics.com/single-cell-gene-expression/datasets/2.1.0/pbmc4k>). Chen et al.'s dataset ([26], GSE87544) of adult mouse hypothalamus cells is available at <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE87544>. Loh et al.'s dataset ([27], SRP 073808) of human embryonic stem cells is available at <https://www.ncbi.nlm.nih.gov/sra/?term=SRP073808>. Petropoulos et al.'s dataset ([28], E-MTAB-3929) of human preimplantation embryos cells is available at <https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-3929>. All the source codes and supplementary data are available at <https://github.com/jianghuaijie/scLRTC>.

Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Received: 28 July 2021 Accepted: 13 October 2021

Published online: 29 November 2021

References

- Chen C, Wu C, Wu L, Wang X, Deng M, Xi R. scRMD: Imputation for single cell RNA-seq data via robust matrix decomposition. *Bioinformatics*. 2020; 36(10):3156–61. <https://doi.org/10.1093/bioinformatics/btaa139>.
- Zhu K, Anastassiou D. 2DImpute: imputation in single-cell RNA-seq data from correlations in two dimensions. *Bioinformatics*. 2020;36(11):3588–9. <https://doi.org/10.1093/bioinformatics/btaa148>.
- Aridakessian C, Poirion O, Yunits B, Zhu X, Garmire LX. DeepImpute: an accurate, fast, and scalable deep neural network method to impute single-cell RNA-seq data. *Genome Biol*. 2019;20(1):1–14. <https://doi.org/10.1186/s1059-019-1837-6>.
- Eraslan G, Simon LM, Mircea M, Mueller NS, Theis FJ. Single-cell RNA-seq denoising using a deep count autoencoder. *Nat Commun*. 2019;10(1):1–14. <https://doi.org/10.1038/s41467-018-07931-2>.
- Xu Y, Zhang Z, You L, Liu J, Fan Z, Zhou X. scGANs: single-cell RNA-seq imputation using generative adversarial networks. *Nucleic Acids Res*. 2020; 48(15):e85. <https://doi.org/10.1093/nar/gkaa506>.
- Andrews TS, Hemberg M. False signals induced by single-cell imputation. *F1000Res*. 2018;7:1740.
- Huang M, Wang J, Torre E, Dueck H, Shaffer S, Bonasio R, et al. SAVER: gene expression recovery for single-cell RNA sequencing. *Nat Methods*. 2018; 15(7):539–42. <https://doi.org/10.1038/s41592-018-0033-z>.
- Dijk D, Sharma R, Nainys J, Yim K, Kathail P, Carr AJ, et al. Recovering gene interactions from single-cell data using data diffusion. *Cell*. 2018;174(3):716–29. <https://doi.org/10.1016/j.cell.2018.05.061>.
- Li WW, Li JJ. An accurate and robust imputation method scImpute for single-cell RNA-seq data. *Nat Commun*. 2018;9(1):1–9.
- Gong W, Kwak IY, Pota P, Koyano-Nakagawa N, Garry DJ. DrImpute: imputing dropout events in single cell RNA sequencing data. *BMC Bioinforma*. 2018;19(1):1–10. <https://doi.org/10.1186/s12859-018-2226-y>.
- Xu J, Cai L, Liao B, Zhu W, Yang J. CMF-impute: an accurate imputation tool for single-cell RNA-seq data. *Bioinformatics*. 2020;36(10):3139–47. <https://doi.org/10.1093/bioinformatics/btaa109>.

12. Linderman GC, Zhao J, Kluger Y. Zero-preserving imputation of scRNA-seq data using low-rank approximation. *bioRxiv*. 2018:397588.
13. Mongia A, Sengupta D, Majumdar A. McImpute: matrix completion based imputation for single cell RNA-seq data. *Front Genet*. 2019;10:9. <https://doi.org/10.3389/fgene.2019.00009>.
14. Ruchansky N, Crovella M, Terzi E. Targeted matrix completion. In: Proceedings of the 2017 SIAM International Conference on Data Mining; 2017. p. 255–263. <https://doi.org/10.1137/1.9781611974973.29>.
15. Zhang L, Zhang S. Imputing single-cell RNA-seq data by considering cell heterogeneity and prior expression of dropouts. *J Mol Cell Biol*. 2021;13(1): 29–40. <https://doi.org/10.1093/jmcb/mjaa052>.
16. Ni Z, Zheng X, Zheng X, Zou X. scLRD: A novel low rank tensor decomposition method for imputing missing values in single-cell multi-omics sequencing data. *IEEE ACM Trans Comput Biol Bioinform*. 2020. p. 1–1. <https://doi.org/10.1109/TCBB.2020.3025804>.
17. Liu J, Musialski P, Wonka P, Ye J. Tensor completion for estimating missing values in visual data. *IEEE Trans Pattern Anal Mach Intell*. 2012;35(1):208–20. <https://doi.org/10.1109/TPAMI.2012.39>.
18. Hu Y, Li B, Zhang W, Liu N, Cai P, Chen F. WEDGE: imputation of gene expression values from single-cell RNA-seq datasets using biased matrix decomposition. *Brief Bioinform*. 2021;22(5). <https://doi.org/10.1093/bib/bbab085>.
19. Wang J, Ma A, Chang Y, Gong J, Jiang Y, Qi R, et al. scGNN is a novel graph neural network framework for single-cell RNA-Seq analyses. *Nat Commun*. 2021;12(1):1882. <https://doi.org/10.1038/s41467-021-22197-x>.
20. Pollen AA, Nowakowski TJ, Shuga J, Wang X, Leyrat AA, Lui JH, et al. Low-coverage single-cell mRNA sequencing reveals cellular heterogeneity and activated signaling pathways in developing cerebral cortex. *Nat Biotechnol*. 2014;32(10):1053–8. <https://doi.org/10.1038/nbt.2967>.
21. Usoskin D, Furlan A, Islam S, Abdo H, Lönnerberg P, Lou D, et al. Unbiased classification of sensory neuron types by large-scale single-cell RNA sequencing. *Nat Neurosci*. 2015;18(1):145–53. <https://doi.org/10.1038/nn.3881>.
22. Yan L, Yang M, Guo H, Yang L, Wu J, Li R, et al. Single-cell RNA-Seq profiling of human preimplantation embryos and embryonic stem cells. *Nat Struct Mol Biol*. 2013;20(9):1131–9. <https://doi.org/10.1038/nsmb.2660>.
23. Zeisel A, Muñoz-Manchado AB, Codeluppi S, Lönnerberg P, La Manno G, Juréus A, et al. Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. *Science*. 2015;347(6226):1138–42. <https://doi.org/10.1126/science.aaa1934>.
24. Han X, Wang R, Zhou Y, Fei L, Sun H, Lai S, et al. Mapping the mouse cell atlas by microwell-seq. *Cell*. 2018;172(5):1091–107. <https://doi.org/10.1016/j.cell.2018.02.001>.
25. Zheng GX, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, et al. Massively parallel digital transcriptional profiling of single cells. *Nat Commun*. 2017;8(1):1–12. <https://doi.org/10.1038/ncomms14049>.
26. Chen R, Wu X, Jiang L, Zhang Y. Single-cell RNA-seq reveals hypothalamic cell diversity. *Cell Rep*. 2017;18(13):3227–41. <https://doi.org/10.1016/j.celrep.2017.03.004>.
27. Loh KM, Chen A, Koh PW, Deng TZ, Sinha R, Tsai JM, et al. Mapping the pairwise choices leading from pluripotency to human bone, heart, and other mesoderm cell types. *Cell*. 2016;166(2):451–67. <https://doi.org/10.1016/j.cell.2016.06.011>.
28. Petropoulos S, Edsgård D, Reinius B, Deng Q, Panula SP, Codeluppi S, et al. Single-cell RNA-seq reveals lineage and X chromosome dynamics in human preimplantation embryos. *Cell*. 2016;165(4):1012–26. <https://doi.org/10.1016/j.cell.2016.03.023>.
29. Zappia L, Phipson B, Oshlack A. Splatter: simulation of single-cell RNA sequencing data. *Genome Biol*. 2017;18(1):1–15. <https://doi.org/10.1186/s13059-017-1305-0>.
30. Kiselev VY, Kirschner K, Schaub MT, Andrews T, Yiu A, Chandra T, et al. SC3: consensus clustering of single-cell RNA-seq data. *Nat Methods*. 2017;14(5): 483–6. <https://doi.org/10.1038/nmeth.4236>.
31. Strehl A, Ghosh J. Cluster ensembles—a knowledge reuse framework for combining multiple partitions. *J Mach Learn Res*. 2002;3(12):583–617.
32. Rand WM. Objective criteria for the evaluation of clustering methods. *J Am Stat Assoc*. 1971;66(336):846–50. <https://doi.org/10.1080/01621459.1971.10482356>.
33. Van der Maaten L, Hinton G. Visualizing data using t-SNE. *J Mach Learn Res*. 2008;9:2579–605.
34. Wolf FA, Angerer P, Theis FJ. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol*. 2018;19(1):1–5. <https://doi.org/10.1186/s13059-017-1382-0>.
35. Blondel VD, Guillaume JL, Lambiotte R, Lefebvre E. Fast unfolding of communities in large networks. *J Stat Mech: Theory Exp*. 2008;2008(10): P10008.
36. Wold S, Esbensen K, Geladi P. Principal component analysis. *Chemometr Intell Lab*. 1987;2(1):37–52. [https://doi.org/10.1016/0169-7439\(87\)80084-9](https://doi.org/10.1016/0169-7439(87)80084-9).
37. Becht E, McInnes L, Healy J, Dutertre C, Kwok I, Ng L, et al. Dimensionality reduction for visualizing single-cell data using UMAP. *Nat Biotechnol*. 2019; 37(1):38–44. <https://doi.org/10.1038/nbt.4314>.
38. Rousseeuw PJ. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *J Comput Appl Math*. 1987;20:53–65. [https://doi.org/10.1016/0377-0427\(87\)90125-7](https://doi.org/10.1016/0377-0427(87)90125-7).
39. Finak G, McDavid A, Yajima M, Deng J, Gersuk V, Shalek AK, et al. MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data. *Genome Biol*. 2015;16(1):1–13. <https://doi.org/10.1186/s13059-015-0844-5>.
40. Ji Z, Ji H. TSCAN: Pseudo-time reconstruction and evaluation in single-cell RNA-seq analysis. *Nucleic Acids Res*. 2016;44(13):e117. <https://doi.org/10.1093/nar/gkw430>.
41. Wang B, Zhu J, Pierson E, Ramazzotti D, Batzoglou S. Visualization and analysis of single-cell RNA-seq data by kernel-based similarity learning. *Nat Methods*. 2017;14(4):414–6. <https://doi.org/10.1038/nmeth.4207>.
42. De Lathauwer L, De Moor B, Vandewalle J. A multilinear singular value decomposition. *SIAM J Matrix Anal Appl*. 2000;21(4):1253–78. <https://doi.org/10.1137/S0895479896305696>.
43. Cai JF, Candès EJ, Shen Z. A singular value thresholding algorithm for matrix completion. *SIAM J Optim*. 2010;20(4):1956–82. <https://doi.org/10.1137/080738970>.
44. Ma S, Goldfarb D, Chen L. Fixed point and Bregman iterative methods for matrix rank minimization. *Math Program*. 2011;128(1):321–53. <https://doi.org/10.1007/s10107-009-0306-5>.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

