1   **A Universal Deep Neural Network for In-Depth Cleaning of Single-Cell RNA-Seq**

2   **Data**

3   Hui Li[1,2], Cory R. Brouwer[1,2], Weijun Luo[1,2*]

4   [1]Department of Bioinformatics and Genomics, College of Computing and Informatics, UNC Charlotte,

5   Charlotte, NC 28223

6   [2]UNC Charlotte Bioinformatics Service Division, North Carolina Research Campus, Kannapolis, NC 28081

7   [*]Correspondence: Weijun.Luo@uncc.edu

8   **Abstract**

9   Single cell RNA sequencing (scRNA-Seq) has been widely used in biomedical research and generated

10   enormous volume and diversity of data. The raw data contain multiple types of noise and

11   technical artifacts and need thorough cleaning. The existing denoising and imputation methods

12   largely focus on a single type of noise (i.e. dropouts) and have strong distribution assumptions

13   which greatly limit their performance and application. We designed and developed the

14   AutoClass model, integrating two deep neural network components, an autoencoder and a

15   classifier, as to maximize both noise removal and signal retention. AutoClass is free of

16   distribution assumptions, hence can effectively clean a wide range of noises and artifacts.

17   AutoClass outperforms the state-of-art methods in multiple types of scRNA-Seq data analyses,

18   including data recovery, differential expression analysis, clustering analysis and batch effect

19   removal. Importantly, AutoClass is robust on key hyperparameter settings including bottleneck

20   layer size, pre-clustering number and classifier weight. We have made AutoClass open source at:

21   https://github.com/datapplab/AutoClass.

22   **Introduction**

23   scRNA-Seq has been widely adopted in biological and medical research[1-5] as an ultra-high

24   resolution and ultra-high throughput transcriptome profiling technology. Enormous amount of

25   data has been generated providing great opportunities and challenges in data analytics.

1   First of all, scRNA-Seq data come with multiple types of noise and quality issues. Some are

2   issues associated with gene expression profiling in general, including RNA amplification bias,

3   uneven library size, sequencing and mapping error, etc. Others are specific to single cell assays.

4   For example, extremely small sample quantity and low RNA capture rate result in large number

5   of false zero expression or dropout[6]. Individual cells vary in differentiation or cell cycle stages[7],

6   health conditions or stochastic transcription activities, which are biological differences but

7   irrelevant in most studies. In addition, substantial batch effects are frequently observed[8] due to

8   inconsistence in sample batches and experiments. Most of these noises and variances are not

9   dropout and may follow Gaussian, Poisson or more complex distributions depending on the

10  source of the variances. All of these variances need to be corrected and cleaned so that

11  biologically relevant differences can be reconstructed and analyzed accurately.

12  Multiple statistical methods have been developed to impute and denoise scRNA-Seq data. Most

13  of these methods rely on distribution assumptions on scRNA-Seq data matrix. For example,

14  deep count autoencoder (DCA)[9] assumes negative binomial distribution with or without zero

15  inflation, SAVER[10] assumes negative binomial distribution, and scImpute[11] uses a mixture of

16  Gaussian and Gamma model. Currently, there is no consensus on the distribution of scRNA-Seq

17  data. Method with inaccurate distribution assumptions[12] may not denoise properly, but rather

18  introduce new complexities and artifacts. Importantly, these methods largely focus on dropouts

19  and ignore other types of noise and variances, which hinders accurate analysis and

20  interpretation of the data.

21  To address these issues, we developed AutoClass, a neural network-based method. AutoClass

22  integrates two neural network components: an autoencoder and a classifier (Figure 1a and

23  Methods). The autoencoder itself consists of two parts: an encoder and a decoder. The encoder

24  reduces data dimension and compresses the input data by decreasing hidden layer size

25  (number of neurons). The decoder, in the opposite, expands data dimension and reconstructs

26  the original input data from the compressed data by increasing hidden layer size. Note the

27  encoder and decoder are symmetric in both architecture and function. The data is most

28  compressed at the so-called bottleneck layer between the encoder and the decoder. The

29  autoencoder itself, as an unsupervised data reduction method, is not sufficient in separating

1    signal from noise (Figure 1b). To ensure the encoding process filter out noise and retain signal,

2    we add a classifier branch from the bottleneck layer (Figure 1a and Methods). When cell classes

3    are unknown, virtual class labels are generated by pre-clustering. Therefore, AutoClass is a

4    composite deep neural network with both unsupervised (autoencoder) and supervised

5    (classifier) learning components. AutoClass does not presume any type of data distribution,

6    hence has the potential to correct a wide range noises and non-signal variances. In addition, it

7    can model non-linear relationships between genes with non-linear activation functions. In this

8    study, we extensively evaluated AutoClass against existing methods using multiple simulated

9    and real datasets. We demonstrated AutoClass can better reconstruct scRNA-Seq data and

10   enhance downstream analysis in multiple aspects. In addition, AutoClass is robust over

11   hyperparameter settings and the default setting applies well in various datasets and conditions.


12   **Results**

13   **Validation of the classifier component**

14   The unique part of AutoClass is the classifier branch from the bottleneck layer. Since encoding

15   process losses information in the input data, the classifier branch is added to make sure

16   relevant information or signal is sufficiently retained. To show that the classifier is needed, we

17   simulated a scRNA-Seq Dataset 1 (see Methods and Supplementary Table 2) using Splatter[13]

18   with 1,000 genes and 500 cells in 6 groups, with and without dropout. Applied both AutoClass

19   and a regular autoencoder without the classifier on the data with dropout, the results are

20   illustrated in two-dimensional t-SNE (see Methods) plots in Figure 1b. AutoClass but not the

21   regular autoencoder was able to recover cell type pattern, indicating the classifier component is

22   necessary for reconstructing scRNA-Seq data.

23   **Gene expression data recovery**

24   We evaluated expression value recovery on simulated scRNA-Seq data with different noise

25   types or distributions. We generated and scRNA-Seq dataset using Splatter with 500 cells, 1000

26   genes in five cell groups with (raw data, Dataset 2) and without dropout (true data). From the

27   same true data, we also generated 5 additional raw datasets by adding noise following different

1   distributions which are representative and commonly seen, including random uniform (Dataset

2   3), Gaussian (Dataset 4), Gamma (Dataset 5), Poisson (Dataset 6) and negative binomial

3   (Dataset 7) (details in Methods and Supplement Table 2 and 3).

4   As expected, dropout noise greatly reduced the data quality and obscured the signal or

5   biological differences such as distinction between cell types (Figure 2a). All other noise types

6   had similar effect on the data (Figure2b-2c and Supplementary Figure 1). With t-SNE

7   transformation on Dataset 2-7, the true data without noise showed distinct cell types, but not

8   the raw data with noises (Figure 2a-c and Supplementary Figure 1). The average Silhouette

9   width[14] (ASW) on the t-SNE plot is a measurement of distance between groups, ranges from -1

10  to 1, where higher values indicate more confident clustering. ASW dropped greatly from 0.64 to

11  around 0 in all raw datasets. After imputation by AutoClass, the cell type pattern was recovered

12  and ASW increased back substantially to 0.2-0.5. In contrast, all published control methods

13  (DCA, MAGIC[15], scImpute and SAVER) were unable to recover the original cell type pattern

14  (Figure 2a-c and Supplementary Figure 1) and ASW scores remained low (Figure 2d) for all noise

15  types.

16  We also measure the data recovery quality using other metrics. The mean squared error (MSE)

17  between the true data and imputed/denoised data for Dataset 3-7 (5 noise types other than

18  dropout) were also computed (Figure 2e). Among the 5 tested methods, AutoClass consistently

19  achieved the smallest MSE for all noise types (Figure 2e). Dropout noise (Dataset 2) is very

20  different from all other noise types (Dataset 3-7) in both distribution form and generation

21  mechanism, and MSE was not an informative measurement of data recovery. We computed the

22  average recovered values of dropout zeros and those of true zeros (Figure 2f) instead. An ideal

23  imputation method can distinguish between these two types of zeros, i.e. impute dropout zeros

24  while retain true zeros (Figure 2f). While SAVER was too conservative in imputing both types of

25  0 values, DCA and MAGIC were too aggressive. AutoClass and scImpute both achieved good

26  balance between imputing dropout 0s and retaining true 0s, yet only the former but bot the

27  later was able to recover the biological difference or distinct cell type clustering (Figure 2a-c

28  and Supplementary Figure1).

1

## Differential expression analysis

Differential expression (DE) analysis is by far the most common analysis of scRNA-Seq and gene expression data. To study the performance of AutoClass in DE analysis, we simulated a scRNA-Seq Dataset 8 using Splatter with 1,000 genes and 500 cells in two cell groups. Here the ground truth of truly differentially expressed genes is known. We applied Two-sample T-test to the true, raw and imputed data using different methods. The median value of t-statistics for the truly differentially expressed genes dropped from 5.79 in the true data to 2.11 in the raw data, and increased back to 5.86 upon imputation by AutoClass, which was almost the same as in the true data and higher than in all control methods (Figure 3 a − b). As shown by ROC curves and area under the curves (AUC), AutoClass also was the best at balancing true positives and false negatives (Figure 3 c − d).

Similarly, AutoClass can improve DE analysis in data with Gaussian noise. We manually added Gaussian noise to the true data of Dataset 8 to generate the raw data of Dataset 9. The DE analysis results can be found in Supplementary Figure 2.

AutoClass also improves marker gene expression analysis. Baron dataset[16] provides known marker gene lists for related cell types in pancreatic islets. AutoClass imputed data increased both t-statistics and fold changes of the marker genes (Figure 3 e − f).

## Clustering analysis

Clustering analysis is frequently done on scRNA-seq data as to identify cell types or subpopulations. To evaluate AutoClass for clustering analysis, we used four real datasets, including two small datasets: the Buettner dataset[2] (182 cells) and the Usoskin dataset[17] (622 cells) and two large datasets: the Lake dataset[18] (8,592 cells) and the Zeisel dataset[19] (3,005 cells). Detailed information for these datasets can be found in Methods and Supplementary Table 1.

1    We compared K-means clustering results on the 200 highest variable genes. The ground truth

2    or the actual number of cell types were used as number of clusters. Clustering results were

3    evaluated by four different metrics: adjusted Rand index[20] (ARI), Jaccard Index[21] (JI), normalized

4    mutual information[22] (NMI) and purity score[23] (PS). All of them range from 0 to 1, with 1

5    indicating a perfect match to the true groups. AutoClass is the only method improving all four

6    metrics from the raw data. In addition, AutoClass achieved the best clustering results for 3 out

7    of 4 datasets (Table 1).

8

| metric | dataset | Raw | AutoClass | DCA | MAGIC | scImpute | SAVER |
|--------|---------|-----|-----------|-----|-------|----------|-------|
| ARI | Buettner | 0.023 | **0.372** | 0.288 | 0.213 | 0.039 | 0.016 |
|     | Usoskin | 0.221 | **0.869** | 0.234 | 0.813 | 0.067 | 0.317 |
|     | Lake | 0.403 | 0.557 | **0.572** | 0.440 | 0.313 | 0.465 |
|     | Zeisel | 0.737 | **0.793** | 0.753 | 0.433 | 0.623 | 0.763 |
| JI | Buettner | 0.242 | **0.409** | 0.363 | 0.368 | 0.262 | 0.247 |
|    | Usoskin | 0.324 | **0.830** | 0.284 | 0.764 | 0.266 | 0.351 |
|    | Lake | 0.323 | 0.439 | **0.453** | 0.346 | 0.254 | 0.364 |
|    | Zeisel | 0.646 | **0.713** | 0.664 | 0.370 | 0.679 | 0.677 |
| NMI | Buettner | 0.035 | **0.395** | 0.333 | 0.335 | 0.075 | 0.038 |
|     | Usoskin | 0.225 | **0.829** | 0.253 | 0.771 | 0.048 | 0.431 |
|     | Lake | 0.611 | 0.667 | **0.676** | 0.601 | 0.500 | 0.642 |
|     | Zeisel | 0.747 | **0.784** | 0.746 | 0.598 | 0.798 | 0.762 |
| PS | Buettner | 0.434 | **0.720** | 0.648 | 0.599 | 0.445 | 0.423 |
|    | Usoskin | 0.545 | **0.937** | 0.579 | 0.913 | 0.416 | 0.682 |
|    | Lake | 0.723 | **0.772** | 0.766 | 0.693 | 0.610 | 0.742 |
|    | Zeisel | 0.894 | **0.917** | 0.880 | 0.763 | 0.548 | 0.897 |

9    **Table 1** Evaluation of clustering results of four real scRNA-Seq datasets. The four metrics are adjusted Rand index (ARI), Jaccard
10   Index (JI), normalized mutual information (NMI) and purity score (PS). Highest value in each row was highlighted in boldface.

11   For the Usoskin dataset, out of all tested methods, only AutoClass and MAGIC reconstructed

12   distinct clusters (Figure 4a). But MAGIC likely generated false positive signals, given that the

13   between-group cell-to-cell correlation are almost the same as within-group correlation, and

14   both are close to 1 (Figure 4b). AutoClass was the only method differentiating within-group vs

15   between-group correlation as informative metrics for signal vs noise (Figure 4b).

**Batch effect removal**

Batch effect rises from different individual cell donors, sample groups, or experiment conditions and can severely affect downstream analysis. We analyzed two real datasets with major batch effect. The Villani dataset[24] sequenced 768 human blood dendritic cells (DC) in 2 batched using Smart-Seq2. The Baron dataset includes 7,162 pancreatic islet cells from 3 healthy individuals.

Similar to Tran et al.[8], we evaluated the performance of batch effect correction as the ability to merge different batches of the same cell type while keeping different cell types separate. We did t-SNE transformation on the data first (Figure 5a and Supplementary Figure 3), then applied the four metrics above mentioned, i.e. ASW, ARI, NMI and PS on both cell types and batches. While cell-type-level metrics measure cell type separation, 1 - batch level metrics measure the merging between batches of same cell type (Figure 5b and Supplementary Figure 4).

In the Villani dataset (Figure 5), the raw data shows clear separation in both cell types and sample batches. After imputation by AutoClass, while cell types remained well separated, the two batches were evenly mixed up within each cell type. In contrast, SAVER failed to reduce the batch effect, while all other methods even aggravated it (Figure 5).

Note that AutoClass corrects the batch effect without knowing the actual number of cell types. Here, we used the default number of clusters in the pre-clustering step, i.e., [8, 9, 10] (see Methods). This is close to the number of spurious groups counting batches (i.e., 8), but far away from the actual number of cell types, or 4. In other words, AutoClass was not misled by the pre-clustering number and correctly recovered the actual cluster number.

In Baron dataset (Supplementary Figure 3 and 4), AutoClass reduced the batch effect and increased cell type separation simultaneously with the default pre-clustering number too. MAGIC dramatically reduced the differences in both batches and cell type. The batch effect correction by other methods were limited.

**Robustness over major hyperparameters**

7

1    AutoClass, as a composite deep neural network, has multiple hyperparameters. Among them,

2    the most important ones are bottleneck layer size, number of pre-clusters and classifier weight.

3    Bottleneck layer plays an important role in autoencoders, it is the narrowest part of the

4    network and the size (number of neurons) controls how much the input data is compressed.

5    The number of clusters ($K$) in the pre-clustering step is specific to the classifier of AutoClass.

6    AutoClass uses three consecutive cluster numbers $[K-1, K, K+1]$, and the final imputation

7    output is the average over three predictions using these three clustering numbers (see

8    Methods). In addition, the classifier weight $w$ (see Equation (4) and Methods) is another

9    AutoClass specific hyperparameter which balance the ratio between autoencoder loss and the

10   classifier loss.

11   AutoClass is robust over a wide range of bottleneck layer sizes, pre-clustering $K$ values (Figure 6

12   and Supplementary Figure 5-6) and classifier weight $w$ (Supplementary Figure 7). The t-SNE

13   clustering patterns, clustering metrics (ASW and ARI), MSE and imputed dropout0s/true 0s ratio

14   remained the same when bottleneck layer size increase from 16 to 256 (Figure 6a, 6c,

15   Supplementary Figure 6a). However, these results or metrics varied heavily in the same analysis

16   using DCA, another autoencoder based method (Figure 6b, 6c, Supplementary Figure 6a).

17   Likewise, AutoClass also achieved stable results over the range of $K$ values – 4-8 (Figure 6d,

18   Supplementary Figure 5, Supplementary Figure 6b) and the range of classifier weight $w$ values –

19   0.1-0.9 (Supplementary Figure 7).


20   **Discussion**

21   In this work, we proposed and developed a deep learning- based method AutoClass for

22   thorough cleaning of scRNA-Seq data. AutoClass integrates two neural network components, an

23   autoencoder and a classifier. This composite network architecture is essential for filtering out

24   noise and retaining signal effectively. Unlike many other scRNA-Seq imputation methods,

25   AutoClass does not rely on any distribution assumption, and fully counts the non-linear

26   interactions between genes. With these properties, AutoClass effectively models and cleans a

27   wide range of noises and artifacts in scRNA-Seq data including dropouts, random uniform,

28   Gaussian, Gamma, Poisson and negative binomial noises, as well as batch effects. These are the

1   most common and representative types of noises and artifacts. Any other types not directly

2   tested would likely be cleaned with the same efficiency because they are similar in distribution

3   and source and AutoClass has no assumption on the noise forms. Such in-depth cleaning led to

4   consistent and substantial improvement of the data quality and downstream analyses including

5   differential expression and clustering, as shown by a range of experiments with both simulated

6   and real datasets.

7   Hyperparameter tuning is an important yet tedious step for training neural network models.

8   Inadequate tuning of hyperparameters may lead to suboptimal results. Remarkably, AutoClass

9   is robust with key hyperparameters including bottleneck layer size ($n$), pre-clustering number

10  ($K$) and classifier weight ($w$). The default setting with $n = 128, K = 9, w = 0.9$ works well for

11  most scRNA-Seq datasets and conditions. This robustness makes AutoClass an appealing method

12  for both performance and practical uses.

13  **Methods**

14  **Architecture of AutoClass**

15  AutoClass integrates two neural network components, an autoencoder and a classifier, to

16  impute scRNA-seq data (Figure 1a). The classifier branch is necessary to preserve signals or

17  biological differences (cell type patterns etc.) from loss in data compression by the encoder.

18  When cell classes are unknown, virtual class labels are generated by pre-clustering using K-

19  means method. The total loss of the entire network is the weighted sum of classifier loss (cross-

20  entropy or CE) and the autoencoder loss (MSE). The activation functions in the hidden layers

21  are all rectified linear unit (ReLU), the activation functions for the output layer of autoencoder

22  and classifier are SoftPlus and SoftMax, respectively.

23  The formulation of AutoClass architecture is:

$$B_k = \text{Encoder}(\bar{X}) \tag{1}$$

$$Y_k = \text{Decoder}(B_k) \tag{2}$$

$$C_k = \text{Classifier}(B_k) \tag{3}$$

9

$$L_k = w \times \text{CE}(\hat{C}_k, C_k) + \tag{4}$$
$$(1 - w) \times \text{MSE}(\bar{X}, Y_k)$$

1   Where $B_k, Y_k, C_k$ and $L_k$ are the bottleneck representation, the output of the decoder hence

2   the autoencoder, the output of the classifier and total loss, respectively. $\hat{C}_k$ is the pre-clustering

3   cell type labels for $k$ clusters. $\bar{X}$ is the input of AutoClass, and has been normalized over library

4   size and followed by a $\log_2$ transformation with pseudo count 1:

$$\bar{X} = \log_2 \left( \text{diag}(s_i)^{-1} X + 1 \right) \tag{5}$$

5   $X$ is the raw count matrix and the size factor $s_i$ for cell $i$ is equal to the library size divided by

6   the median library size across cells. Library size is defined as the total number of counts per cell.

7   The final imputed data is the average prediction of the autoencoder over different cluster

8   numbers in pre-clustering:

$$Y = \text{E}(Y_k|k) \tag{6}$$

9    For all datasets in this manuscript, we used 3 consecutive cluster numbers, or  $k =$

10   $[K - 1, K, K + 1]$, the default value is $K = 9$. The final imputation result was the average

11   results over different $K$s.

12   **AutoClass Implementation and hyperparameter settings**

13   AutoClass is implemented in Python 3 with Keras. *Adam* is used for optimizer with default

14   learning rate 0.001. Learning rate is multiplied by 0.1 if validation loss does not improve for 15

15   epochs. The training stops if there is no improvement for 30 epochs.

16   Although AutoClass works well for small bottleneck layer sizes ($n = 16$, 32 or similar), we set

17   the default value to be $n = 128$, as be conservative and to avoid potential information loss in

18   data compression. This default value was used in all datasets in this paper.

19   AutoClass is stable over different choices of $K$ in pre-clustering as long as $K$ is not extremely far

20   away from the true number of cell clusters. The default value $K = 9$ was used in all datasets in

21   this paper except simulated Dataset 8 and 9, since the true number of cell clusters in these two

22   datasets is 2 which is far smaller than default value 9. Hyperparameter $K$ can be chosen based

23   on prior knowledge of the data or statistical methods like elbow method[25] and Silhouette

1 method[14]. $K$ used in Dataset 8 and 9 was the average of estimations by elbow method and

2 Silhouette method.

3 AutoClass is stable on classifier weight $w$ in the range of 0.1-0.9 (supplementary Figure 7). We

4 found that in general classification loss is far smaller than reconstruction loss (Supplementary

5 Figure 8), to have a better balance between those two losses, we set the default value to be

6 $w = 0.9$. This default value was used in all the datasets in this paper.

7 In addition, overfitting is a common problem in neural network models. Dropout of neurons[26] in

8 the bottleneck layer is used in AutoClass to prevent overfitting. Interestingly, a relatively high

9 dropout rate in AutoClass also helps to correct batch effect. In the batch effect removal

10 analyses, we set dropout rate to be 0.5 in AutoClass, and to be fair, also in DCA. But DCA was

11 unable to remove batch effect (Figure 5, Supplementary Figure 3-4). The default dropout rate

12 0.1 in AutoClass was used in all the other datasets and analyses in this paper.

13 AutoClass hyperparameter settings for all the datasets can be found in Supplementary Table 1-

14 3.

15 **Analysis details**

16 _Noise types other than dropout:_ Dataset 3-7 and Dataset 9 were generated by manually adding

17 noise to the true data of Dataset 2 and Dataset 8, respectively. The noise was first generated by

18 _Python numpy.random_ package with different noise distributions (details in Supplementary

19 Table 3) , and then centered (so that noise mean is 0). The noise was then added to true data,

20 all values were rounded to be integers and negative values set to 0, since scRNA-Seq data raw

21 counts are positive integers.

22 _Highly variable genes_: The highly variable genes in each dataset are ranked by the ratio

23 between gene-wise variance vs mean computed from non-zero values.

24 _t-Distributed stochastic neighbor embedding (t-SNE):_ We applied t-SNE[27] to visualize

25 datasets. We first reduce the number of data dimensions by using the top 50 principle

26 components, and then use _TSNE_ function in the _sklearn.manifold_ package with default settings

27 to further reduce the dimension to 2 for visualization.

1    *Batch effect removal score:* Four clustering metrics ASW, ARI, NMI and PS were used to

2    measure the performance of batch effect correction. We applied ASW to the t-SNE transformed

3    data, and batch effect removal was scored by both cell-type-wise ASW vs 1 – batch-wise ASW

4    (Figure 5b and Supplementary Figure 4). Higher values in both dimensions together denote

5    better batch effect removal. ARI, NMI and PS metrics are used and plotted in the same fashion

6    as ASW. To compute ARI, NMI and PS, K-means clustering was performed first to obtain cluster

7    labels, which were then compared to batch labels and cell type labels. The batch indices were

8    computed for each individual cell type first, and take weighted sum across cell types. The

9    weight for each cell type is proportional to the number of cells.

10    **Control methods**

11    DCA (version 0.2) was downloaded from https://github.com/theislab/dca

12    Magic (version 0.1.0) was downloaded from https://github.com/KrishnaswamyLab/MAGIC

13    scImpute (version 0.0.5) was downloaded from https://github.com/Vivianstats/scImpute.

14    SAVER (version 0.3.0) was downloaded from https://github.com/mohuangx/SAVER.

15    **Real scRNA-seq datasets**

16    We collected and analyzed multiple real scRNA datasets from published studies. These datasets

17    have been well established, widely used and tested as shown in literature. While major

18    technical attributes are summarized in Supplementary Table 1, below are more details.

19    *Baron Study*

20    Human pancreatic islets cells data were obtained from 3 healthy individuals, which provided

21    gene expression profiles for 17,434 genes in 7,729 cells. We filtered out genes expressed in less

22    than 5 cells, removed cell types less than 1% of the cell population. Analysis was restricted to

23    top 1,000 highly variable genes. Final dataset contained 7,162 cells with 8 different cell types.

24    The raw counts data are available at

25    https://shenorrlab.github.io/bseqsc/vignettes/pages/data.html.

12

1 *Villani Study*

2 The human blood dendritic data contained 26,593 genes in 1,140 cells. We kept batch 1 (plate

3 id: P10, P7, P8 and P9) batch 2 (plate id P3, P4, P13, P14) cells, and filtered out genes expressed

4 in less than 5 cells. Analysis was restricted to top 1,000 highly variable genes. Final dataset

5 contained 768 cells with 4 different cell types in 2 batches. The raw data are available at GEO

6 accession GSE80171.

7 *Lake Study*

8 Human brain frontal cortex data contained 34,305 genes in 10,319 cells. We filtered out genes

9 expressed in less than 5 cells, removed cell types h less than 3% of the cell population. Analysis

10 was restricted to top 1,000 highly variable genes. Final dataset contained 8,592 cells with 11

11 different cell types. The raw data are available at GEO accession GSE97930.

12 *Zeisel Study*

13 Mouse cortex and hippocampus data contained 19,972 genes in 3,005 cells. We filtered out

14 genes expressed in less than 5 cells. Analysis was restricted to top 1,000 highly variable genes.

15 Final dataset contained 3,005 cells with 9 different cell types. Annotated data are available at

16 http://linnarssonlab.org/cortex.

17 *Buettner Study:*

18 Mouse embryonic stem cells contained 8,989 genes in 182 cells. We filtered out genes

19 expressed in less than 5 cells. Final dataset contained 8,985 genes and 182 cells in 3 cells lines.

20 The full dataset was deposited at ArrayExpress: E-MTAB-2805. The normalized data can be

21 obtained from https://www.nature.com/articles/nbt.3102.

22 *Usoskin Study*

23 Neuronal data contained 17,772 genes in 622 cells. We filtered out genes expressed in less than

24 5 cells. Analysis was restricted to top 1,000 highly variable genes. Final dataset contains 622

25 cells with 4 different cell types. The normalized data can be obtained from

26 https://www.nature.com/articles/nbt.3102.

13

**Simulated scRNA-seq datasets**

Splatter R (version v1.2.2) package was used to simulate scRNA-seq datasets with dropout values. Gaussian noise was manually added when needed. Genes expressed in less than 3 cells were filtered out before analysis. The parameter settings for simulation are summarized in Supplementary Table 2 and 3.

**Code availability**

AutoClass python module, documentation, tutorial with example, and code to reproduce the main results in the manuscript are available online: https://github.com/datapplab/AutoClass.

**Data availability**

All simulated datasets can be generated using the parameters specified in the Simulated scRNA-Seq datasets subsection, all the real datasets are publicly available as mentioned in the Real scRNA-Seq datasets subsection. In addition, multiple simulated and real datasets were provided in the GitHub repository above as demo datasets, ready for analysis.

**References**

1.  Griffiths, J.A., Scialdone, A. & Marioni, J.C. Using single-cell genomics to understand developmental processes and cell fate decisions. *Mol Syst Biol* **14**, e8046 (2018).
2.  Buettner, F. et al. Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations of cells. *Nat Biotechnol* **33**, 155-160 (2015).
3.  Papalexi, E. & Satija, R. Single-cell RNA sequencing to explore immune cell heterogeneity. *Nat Rev Immunol* **18**, 35-45 (2018).
4.  Keren-Shaul, H. et al. A unique microglia type associated with restricting development of Alzheimer's disease. *Cell* **169**, 1276-1290 e1217 (2017).
5.  Stubbington, M.J.T., Rozenblatt-Rosen, O., Regev, A. & Teichmann, S.A. Single-cell transcriptomics to explore the immune system in health and disease. *Science* **358**, 58-63 (2017).
6.  Hicks, S.C., Townes, F.W., Teng, M. & Irizarry, R.A. Missing data and technical variability in single-cell RNA-sequencing experiments. *Biostatistics* **19**, 562-578 (2018).
7.  Trapnell, C. et al. The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat Biotechnol* **32**, 381-386 (2014).
8.  Tran, H.T.N. et al. A benchmark of batch-effect correction methods for single-cell RNA sequencing data. *Genome Biol* **21**, 12 (2020).
9.  Eraslan, G., Simon, L.M., Mircea, M., Mueller, N.S. & Theis, F.J. Single-cell RNA-seq denoising using a deep count autoencoder. *Nat Commun* **10**, 390 (2019).

1   10.   Huang, M. et al. SAVER: gene expression recovery for single-cell RNA sequencing. *Nat Methods*
2         **15**, 539-542 (2018).
3   11.   Li, W.V. & Li, J.J. An accurate and robust imputation method scImpute for single-cell RNA-seq
4         data. *Nat Commun* **9**, 997 (2018).
5   12.   Svensson, V. Droplet scRNA-seq is not zero-inflated. *Nat Biotechnol* **38**, 147-150 (2020).
6   13.   Zappia, L., Phipson, B. & Oshlack, A. Splatter: simulation of single-cell RNA sequencing data.
7         *Genome Biol* **18**, 174 (2017).
8   14.   Rousseeuw, P.J. Silhouettes: a graphical aid to the interpretation and validation of cluster
9         analysis *Journal of Computational and Applied Mathematics* **20**, 53-65 (1987).
10  15.   Dijk, D.v. et al. Recovering gene interactions from single-cell data using data diffusion. *Cell* **174**,
11        716-729 (2017).
12  16.   Baron, M. et al. A single-cell transcriptomic map of the human and mouse pancreas reveals
13        Inter- and Intra-cell population structure. *Cell Syst* **3**, 346-360 e344 (2016).
14  17.   Usoskin, D. et al. Unbiased classification of sensory neuron types by large-scale single-cell RNA
15        sequencing. *Nat Neurosci* **18**, 145-153 (2015).
16  18.   Lake, B.B. et al. Integrative single-cell analysis of transcriptional and epigenetic states in the
17        human adult brain. *Nat Biotechnol* **36**, 70-80 (2018).
18  19.   Zeisel, A. et al. Brain structure. Cell types in the mouse cortex and hippocampus revealed by
19        single-cell RNA-seq. *Science* **347**, 1138-1142 (2015).
20  20.   Hubert, L. & Arabie, P. Comparing partitions. *Journal of classification* **2**, 193-218 (1985).
21  21.   Leydesdorff, L. On the normalization and visualization of author co-citation data: Salton's Cosine
22        versus the Jaccard index. *Journal of the American Society for Information Science and*
23        *Technology* **59**, 77-85 (2007).
24  22.   Estevez, P.A., Tesmer, M., Perez, C.A. & Zurada, J.M. Normalized mutual information feature
25        selection. *IEEE Trans Neural Netw* **20**, 189-201 (2009).
26  23.   Manning, C.D., Raghavan, P. & Schutze, H. Introduction To Information Retrieval. (Cambridge
27        Univ. Press, Cambridge; 2008).
28  24.   Villani, A.C. et al. Single-cell RNA-seq reveals new types of human blood dendritic cells,
29        monocytes, and progenitors. *Science* **356** (2017).
30  25.   Ketchen, D.J. & Shook, C.L. The application of cluster Analysis in strategic management research:
31        an analysis and critique. *Strategic Management* **17**, 441-458 (1996).
32  26.   Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout: a simple way
33        to prevent neural networks from overfitting. *J Mach Learn Res* **15** (2014).
34  27.   Maaten, L.v.d. & Hinton, G. Visualizing data using t-SNE. *Journal of machine learning research* **9**,
35        2579-2605 (2008).
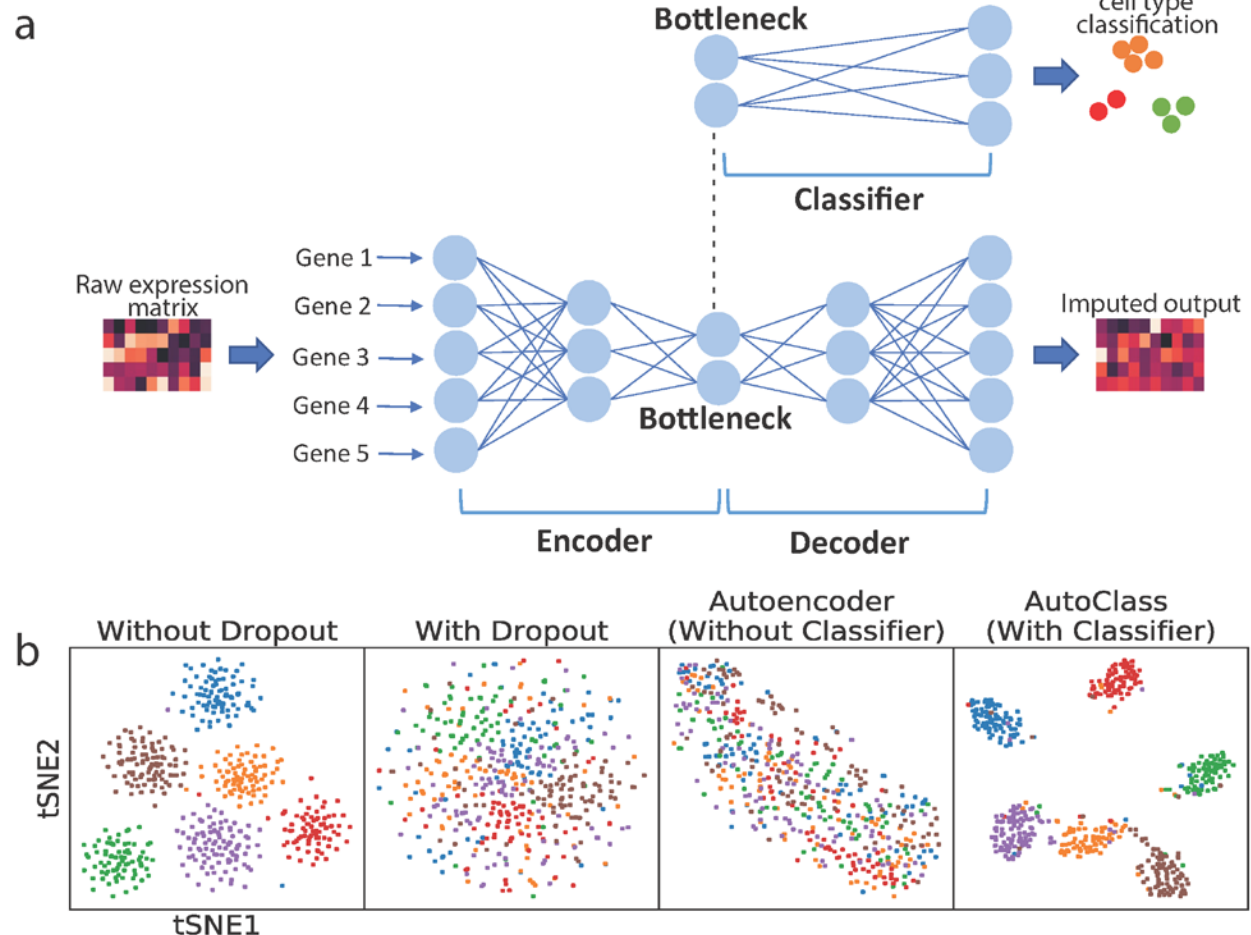
36

37

1  **Figures**

2



3

4  **Figure 1** AutoClass integrates a classifier to a regular autoencoder, as to fully reconstruct scRNA-Seq data. **a** AutoClass consists a
5  regular autoencoder and a classifier branch from the bottleneck layer. The input raw expression data is compressed in the
6  encoder, and reconstructed in the decoder, the classifier branch helps to retain signal in data compression. The output of the
7  autoencoder is the desired imputed data (see Methods for details). **b** t-SNE plots of Dataset 1 without dropout, with dropout,
8  with dropout imputed by a regular autoencoder and AutoClass.

9

**Figure 2** Gene expression data recovery after imputation. **a**, **b** and **c** t-SNE plots for Dataset 2 (dropout noise), Dataset 4 (Gaussian noise) and Dataset 7 (negative binomial noise), respectively. **d** Average Silhouette width based on t-SNE plot for Dataset 2-7. **e** Mean squared error between true data and imputed data for Dataset 3-7. **f** Average recovered values of dropout 0s and true 0s for different imputation methods.
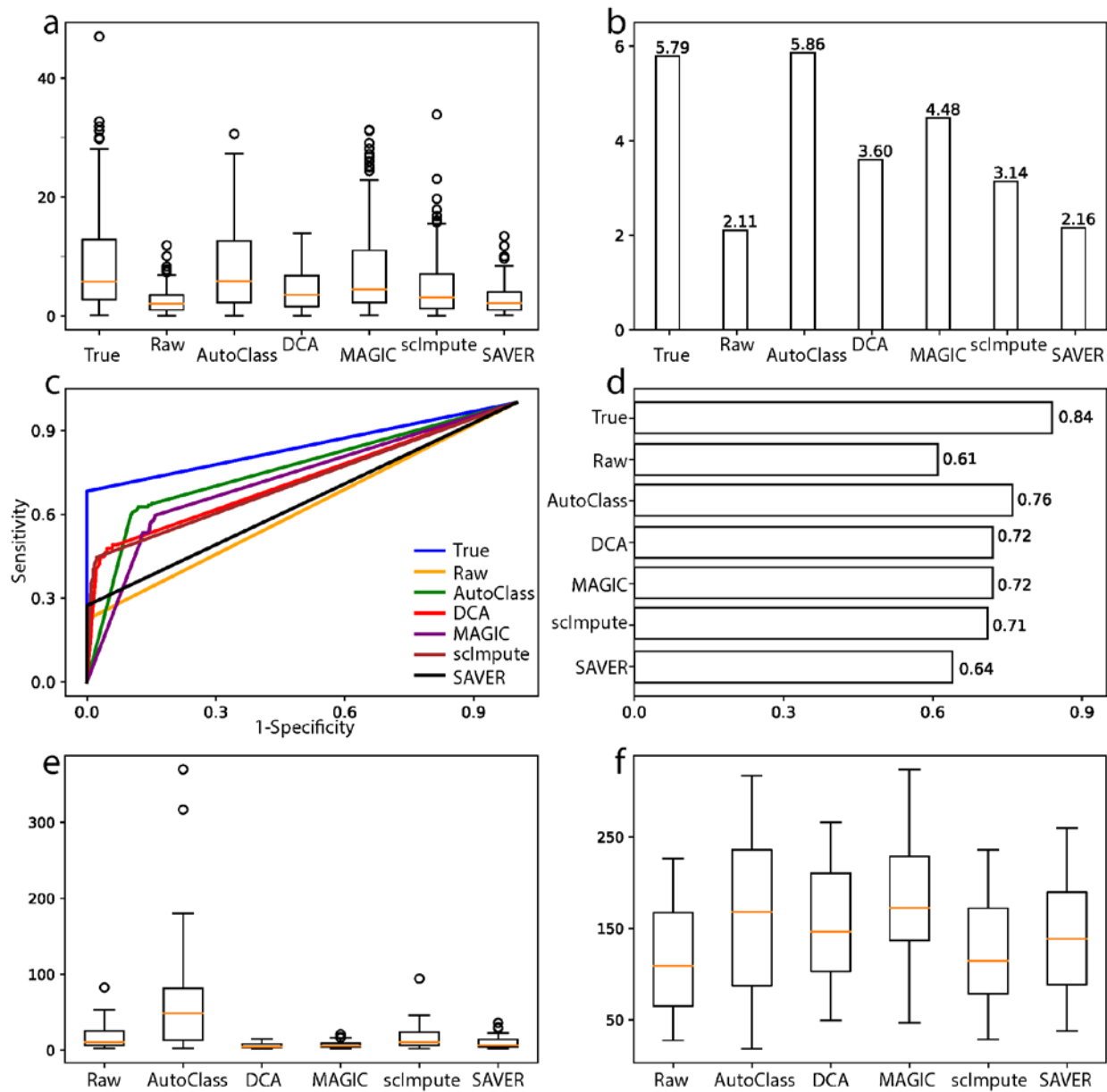
**Figure 3** Differential expression analysis and marker gene analysis. **a** and **b** T-statistics and their median values for truly differentially expressed genes in Dataset 8. **c** and **d** ROC curves and areas under the ROC curves for Dataset 8. **e** and **f** Fold changes and T-statistics of marker genes in the Baron dataset.
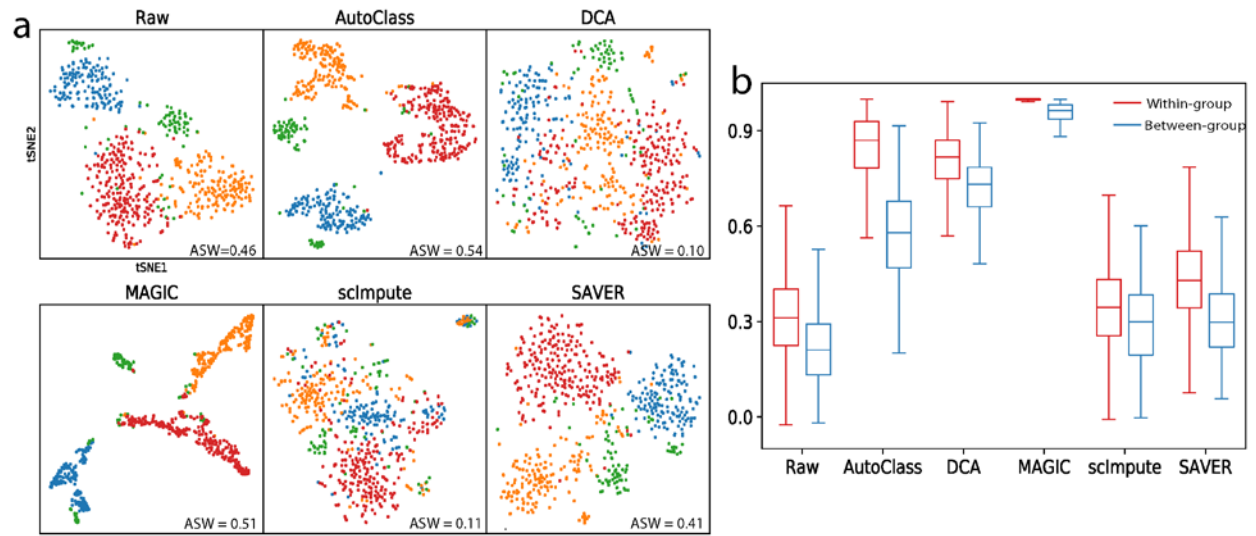
**Figure 4** Imputation results for Usoskin data. **a** t-SNE plots for raw and imputed data. **b** Within-group and between-group cell-to-cell correlation for raw and imputed data.
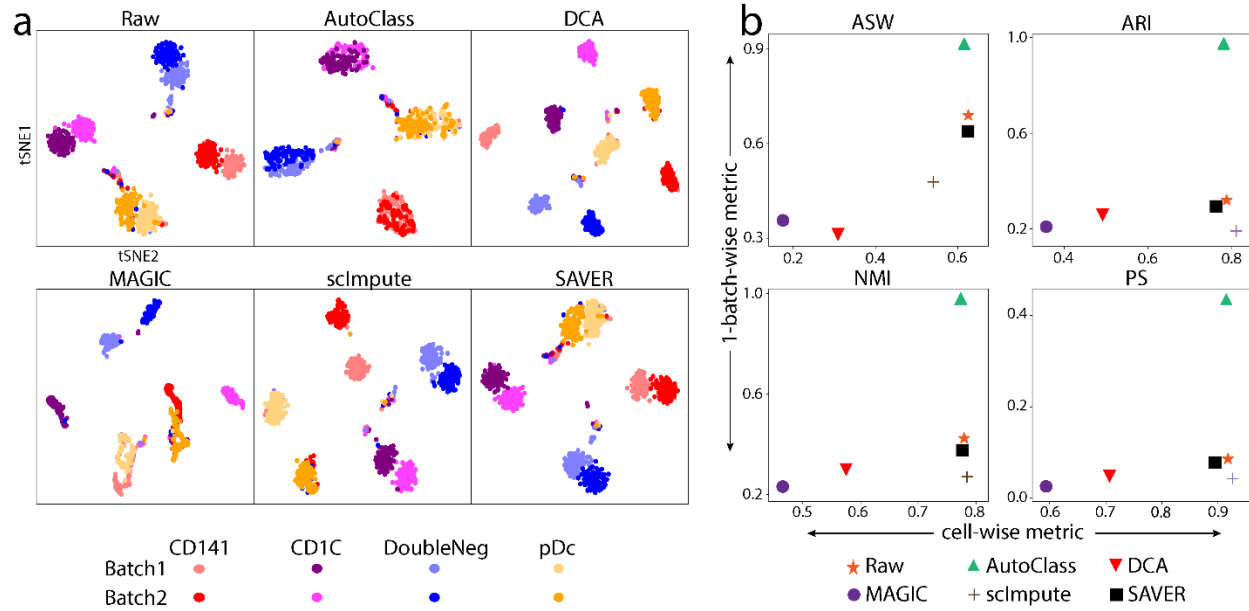
**Figure 5** Batch effect in Villani dataset. **a** t-SNE plots for raw and imputed data. **b** Evaluation of batch and cell type separation in raw and imputed data by four different metrics: average Silhouette width, adjusted Rand index, normalized mutual information and purity score. Good performance is show as big values in both X-(cell type separation) and Y-(batch effect removal) axes.
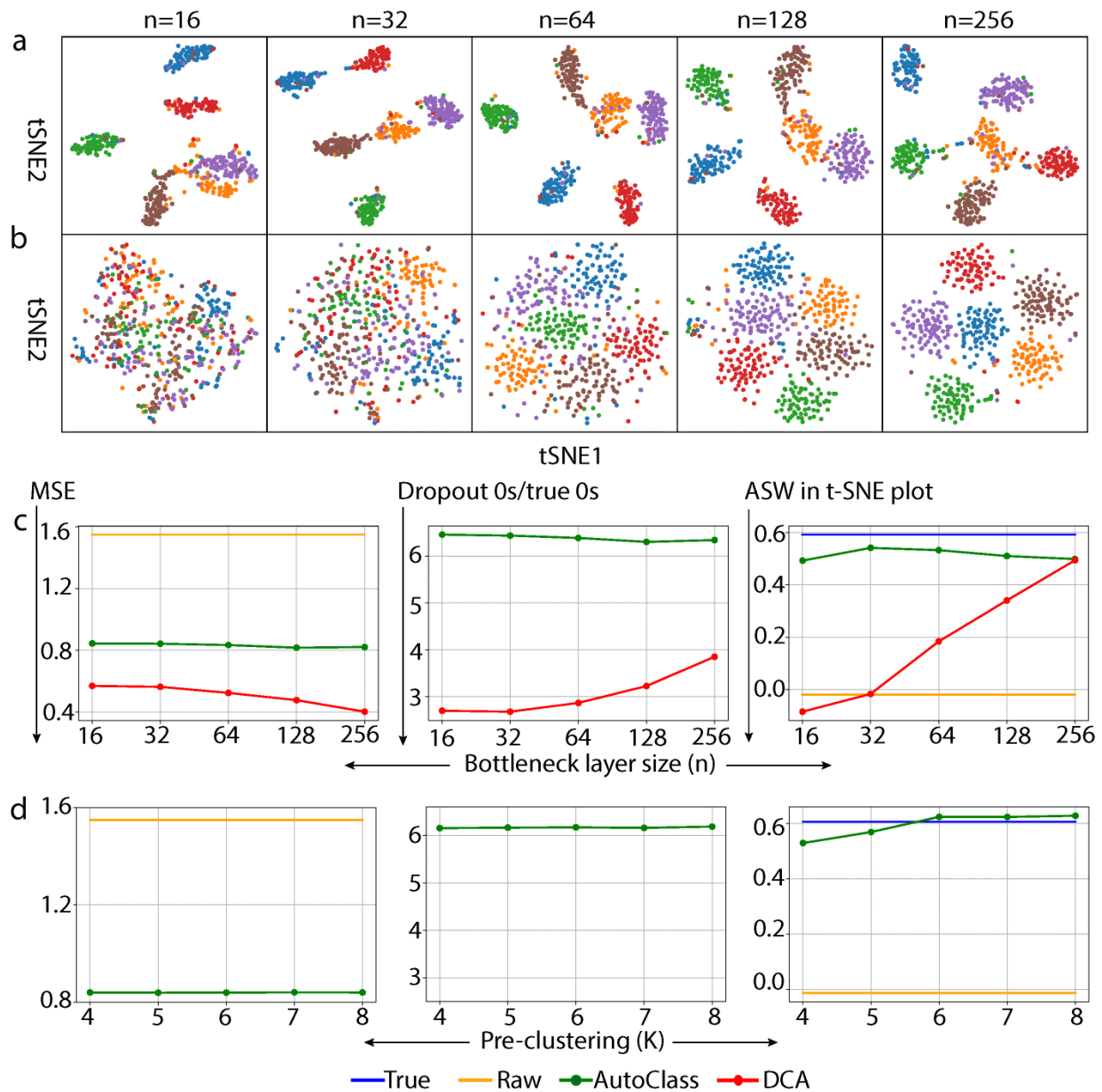
**Figure 6** Impact of bottleneck layer size and the number of clusters in the pre-clustering in Dataset 1. **a** and **b** t-SNE plots of data imputed by AutoClass (**a**) and DCA (**b**) with different bottleneck sizes (n). **c** and **d** Imputation results by AutoClass with bottleneck layer size (**c**) and the number of clusters in the pre-clustering (**d**).