

Data and text mining

scAMACE: model-based approach to the joint analysis of single-cell data on chromatin accessibility, gene expression and methylation

Jiaxuan Wangwu , Zexuan Sun and Zhixiang Lin*

Department of Statistics, The Chinese University of Hong Kong, Hong Kong SAR, China

*To whom correspondence should be addressed.

Associate Editor: Jonathan Wren

Received on March 29, 2021; revised on May 26, 2021; editorial decision on May 31, 2021; accepted on June 3, 2021

Abstract

Motivation: The advancement in technologies and the growth of available single-cell datasets motivate integrative analysis of multiple single-cell genomic datasets. Integrative analysis of multimodal single-cell datasets combines complementary information offered by single-omic datasets and can offer deeper insights on complex biological process. Clustering methods that identify the unknown cell types are among the first few steps in the analysis of single-cell datasets, and they are important for downstream analysis built upon the identified cell types.

Results: We propose scAMACE for the integrative analysis and clustering of single-cell data on chromatin accessibility, gene expression and methylation. We demonstrate that cell types are better identified and characterized through analyzing the three data types jointly. We develop an efficient Expectation–Maximization algorithm to perform statistical inference, and evaluate our methods on both simulation study and real data applications. We also provide the GPU implementation of scAMACE, making it scalable to large datasets.

Availability and implementation: The software and datasets are available at https://github.com/cuhklinlab/scAMACE_py (python implementation) and <https://github.com/cuhklinlab/scAMACE> (R implementation).

Contact: zhixianglin@cuhk.edu.hk

Supplementary information: [Supplementary data](#) are available at *Bioinformatics* online.

1 Introduction

Recent developments in single-cell technologies enable multiple measurements of different genomic features (Lahnemann *et al.*, 2020). Sequencing technologies include single-cell RNA sequencing (scRNA-seq) which measures transcription, single-cell ATAC sequencing (scATAC-seq) and the assay based on combinatorial indexing (sci-ATAC-seq) (Cusanovich *et al.*, 2018b) that measure chromatin accessibility, and single-nucleus methylcytosine sequencing (snmC-seq) (Luo *et al.*, 2017) which measures methylome at the single-cell resolution. High technical variation is presented in single-cell datasets due to the limited amount of genomic materials and the experimental procedures to amplify the signals (Lahnemann *et al.*, 2020).

Because cell types are usually unknown beforehand, clustering methods are needed to identify the cell types. Majority of existing clustering algorithms only take one single dataset as input. Beside the widely used K-Means clustering algorithm, hierarchical clustering (Ward, 1963) forms hierarchical groups of mutually exclusive subsets on the basis of their similarity with respect to specified characteristics by considering the union of all possible $\frac{k(k-1)}{2}$ pairs and

accepting the union with which an optimal value of the objective function is associated. Spectral Clustering (Ng *et al.*, 2001) uses the top k eigenvectors of a matrix derived from the distance between points simultaneously for clustering. Several algorithms are developed specifically for scRNA-seq data. SC3 (Kiselev *et al.*, 2017) combines multiple clustering outcomes through a consensus approach. SIMLR (Wang *et al.*, 2017) learns a distance metric by multiple kernels and clusters with affinity propagation. CIDR (Lin *et al.*, 2017) imputes the gene expression profiles, calculates the dissimilarity based on the imputed gene expression profiles for every pair of single cells, performs principal coordinate analysis using the dissimilarity matrix, and finally performs clustering using the first few principal coordinates. SOUP (Zhu *et al.*, 2019) semi-softly classifies both pure and intermediate cell types: it first identifies the set of pure cells by special block structure and estimates a membership matrix, then estimates soft membership for the other cells. For the analysis of single-cell chromatin accessibility data, scABC (Zamanighomi *et al.*, 2018) first weights cells and applies weighted K-medoids clustering, then calculate landmarks for each cluster, and finally clusters the cells by assignment to the closest landmark based on Spearman correlation. Cusanovich (Cusanovich *et al.*, 2018a)

makes use of singular value decomposition on TF-IDF transformed matrix and density peak clustering algorithm. cisTopic (Bravo González-Blas *et al.*, 2019) uses latent Dirichlet allocation with a collapsed Gibbs sampler to iteratively optimize the region-topic distribution and the topic-cell distribution. SCALE (Xiong *et al.*, 2019) combines the variational autoencoder framework with the Gaussian Mixture Model which extracts latent features that characterize the distributions of input scATAC-seq data, and then uses the latent features to cluster cell mixtures into subpopulations. Clustering methods are also developed for single-cell methylation data. BPRMeth (Kapourani and Sanguinetti, 2016) uses probabilistic machine learning to extract higher order features across a defined region and to cluster promoter-proximal regions by Binomial distributed probit regression (BPR) and mixture modeling. PDclust (Hui *et al.*, 2018) leverages the methylation state of individual CpGs to obtain pairwise dissimilarity (PD) values, and calculates Euclidean distances between each pair of cells using their PD values and performed hierarchical clustering. Melissa (Kapourani and Sanguinetti, 2019) implements a Bayesian hierarchical model that jointly learns the methylation profiles of genomic regions of interest and clusters cells based on their genome-wide methylation patterns. pCSM (Yin *et al.*, 2019) implements a semi-reference-free procedure to perform virtual methylome dissection using the non-negative matrix factorization algorithm. It first determines putative cell-type-specific methylated loci and then clusters the loci into groups based on their correlations in methylation profiles.

Studies based on single-omic data provide only a partial landscape of the entire cellular heterogeneity (Ma *et al.*, 2020). High technical noise and the growth of available datasets measuring different genomic features encourage integrative analysis (Lahnemann *et al.*, 2020). By combining complementary information from multiple datasets, the cell types may be better separated and characterized (Corces *et al.*, 2016; Duren *et al.*, 2017). The integrative analysis of gene expression and chromatin activity may better define cell types and lineages, especially in complex tissues (Duren *et al.*, 2018). Seurat V3 (Stuart *et al.*, 2019) uses Canonical Correlation Analysis (CCA) to reduce the dimension of the datasets. It identifies the pairwise correspondences of single cells across datasets, termed ‘anchors’, and then transfers labels from a reference dataset onto a query dataset. coupleNMF (Duren *et al.*, 2018) is based on the coupling of two non-negative matrix factorizations, where a ‘soft’ clustering can be obtained following the matrix factorizations. It enables integrative analysis of scRNA-seq and scATAC-seq data. LIGER (Welch *et al.*, 2019) integrates multimodal datasets via integrative non-negative matrix factorization (iNMF) to learn a low-dimensional space defined by dataset-specific factors and shared factors across datasets, and then build a neighborhood graph based on the shared factors to identify joint clusters by performing community detection on this graph. scACE (Lin *et al.*, 2020) is a model-based approach that jointly analyzes single-cell chromatin accessibility and scRNA-Seq data, and it quantifies the uncertainty of cluster assignments. MAESTRO (Wang *et al.*, 2020) integrates scRNA-seq and scATAC-seq data from multiple platforms. It also provides comprehensive functions for pre-processing, alignment, quality control and quantification of expression and accessibility. coupleCoC (Zeng *et al.*, 2020) performs co-clustering of the cells and the features simultaneously in the source data and the target data, and it also matches the cell clusters between the source data and the target data through minimizing the distribution divergence. scMC (Zhang and Nie, 2021) integrates multiple scRNA-Seq datasets or multiple scATAC-Seq datasets, where it learns biological variation via variance analysis to subtract technical variation inferred in an unsupervised manner. The three data types, including gene expression, chromatin accessibility and methylation, have distinct characteristics and complex relationships with each other. The aforementioned methods for integrative analysis are not designed to integrate all three data types. Moreover, these methods (except scACE) do not provide statistical inference on the cluster assignments, which may be important when there are cells at the intermediate stages during development.

In this work, we extend scACE (Lin *et al.*, 2020) to scAMACE (integrative Analysis of single-cell Methylation, chromatin Accessibility and gene Expression). scAMACE considers the biological and technical variabilities when integrating multiple data types, and it can provide statistical inference on the assignment of clusters. We reason that by combining complementary biological information from multiple data types, better cell type separation can be achieved. We present our model in Section 2, and statistical inference using the Expectation–Maximization (EM) algorithm in Section 3. Simulation study and real data applications are presented in Sections 4 and 5, respectively. The conclusion is presented in Section 6.

2 Materials and methods

An overview of scAMACE is presented in Figure 1.

2.1 Model for scRNA-Seq data

The model specification for scRNA-Seq data is as the following.

$$\begin{aligned} \omega_{kg}^{rna} \xrightarrow{z_{lk}} u_{lg} &\rightarrow v_{lg} \rightarrow y_{lg} \quad \forall g, \\ P(z_{lk} = 1) &= \psi_k^{rna}, \\ u_{lg} | z_{lk} = 1 &\sim \text{Bernoulli}(\omega_{kg}^{rna}), \\ v_{lg} | u_{lg} &\sim u_{lg} \text{Bernoulli}(\pi_{l1}) + (1 - u_{lg}) \text{Bernoulli}(\pi_{l0}), \\ \pi_{l0} &\sim \text{Beta}(\alpha = 1, \beta = 1), \pi_{l1} \sim \mathbb{1}(\pi_{l1} \geq \pi_{l0}) \text{Beta}(\alpha = 1, \beta = 1), \\ p(y_{lg} | v_{lg}) &= v_{lg} g_1(y_{lg}) + (1 - v_{lg}) g_0(y_{lg}). \end{aligned}$$

We assume that there are K cell clusters in total, the random variable z_{lk} denotes whether cell l belongs to cluster $k \in \{1, \dots, K\}$ and z_l follows categorical distribution with probability ψ_k^{rna} for cluster k .

ω_{kg}^{rna} denotes the probability that gene g is active in cluster k . u_{lg} is a binary latent variable representing whether gene g is active in cell l and $u_{lg} = 1$ represents that it is active. v_{lg} denotes whether gene g is expressed in cell l and $v_{lg} = 1$ represents that it is expressed.

When gene g is active in cell l ($u_{lg} = 1$), the probability that gene g is expressed in cell l ($v_{lg} = 1$) is π_{l1} , while the probability that gene g is expressed is π_{l0} if the gene is not active ($u_{lg} = 0$). Since genes are more likely to be expressed when they are active, we assume that $\pi_{l1} \geq \pi_{l0}$ and the prior distributions of π_{l1} and π_{l0} are assumed to be flat.

Let y_{lg} denote the observed gene expression for gene g in cell l (after normalization to account for sequencing depth and gene length), and we assume that $y_{lg} | v_{lg}$ follows a mixture distribution,

$$\begin{aligned} \omega_{kg}^{acc} \xrightarrow{z_{lk}} u_{ig} &\rightarrow o_{ig} \rightarrow x_{ig} \quad (\text{scCAS data}) \\ \omega_{kg}^{rna} \xrightarrow{z_{lk}} u_{lg} &\rightarrow v_{lg} \rightarrow y_{lg} \quad (\text{scRNA-Seq data}) \\ \omega_{kg}^{met} \xrightarrow{z_{lk}} u_{dg} &\rightarrow m_{dg} \rightarrow t_{dg} \quad (\text{sc-methylation data}) \\ \omega_{kg}^{acc} &\sim \text{Beta}(\mu_{kg}^{acc}, \phi^{acc}), \text{logit}(\mu_{kg}^{acc}) = f(\omega_{kg}^{rna}), \forall g \\ \omega_{kg}^{met} &\sim \text{Beta}(\mu_{kg}^{met}, \phi^{met}), \text{logit}(\mu_{kg}^{met}) = g(\omega_{kg}^{rna}), \forall g \end{aligned}$$

Notations:

- ω : cluster-specific regulatory region / gene activity
- z : cluster assignment
- u (Binary): gene activity status (active or not)
- v (Binary): gene expression status (expressed or not)
- o (Binary): gene activity score level (high or low)
- m (Binary): gene methylated status (methylated or not)
- x : observed gene activity score
- y : observed gene expression level
- t : observed gene methylation level

Fig. 1. Graphical representation of scAMACE

where $g_1(\cdot)$ and $g_0(\cdot)$ are density functions of the expression level conditional on v_{ig} .

2.2 Model for single-cell chromatin accessibility (scCAS) data

The model specification for scCAS data is as the following.

$$\begin{aligned} \omega_{kg}^{acc} &\xrightarrow{z_{ik}} u_{ig} \rightarrow o_{ig} \rightarrow x_{ig} \quad \forall g, \\ P(z_{ik} = 1) &= \psi_k^{acc}, \\ u_{ig}|z_{ik} = 1 &\sim \text{Bernoulli}(\omega_{kg}^{acc}), \\ o_{ig}|u_{ig} &\sim u_{ig}\text{Bernoulli}(\pi_{i1}) + (1 - u_{ig})\text{Bernoulli}(\pi_{i0}), \\ \pi_{i1} &\sim \text{Beta}(\alpha_{acc} = 1, \beta_{acc} = 1), \text{ set } \pi_{i0} = 0, \\ p(x_{ig}|o_{ig}) &= o_{ig}f_1(x_{ig}) + (1 - o_{ig})f_0(x_{ig}), \\ \omega_{kg}^{acc}|\omega_{kg}^{rna} &\sim \text{Beta}(\mu_{kg}^{acc}, \phi^{acc}), \text{ logit}(\mu_{kg}^{acc}) = f(\omega_{kg}^{rna}). \end{aligned}$$

The random variables ω_{kg}^{acc} , z_{ik} , ψ_k^{acc} and u_{ig} have similar interpretations to their corresponding variables in the model for scRNA-Seq data. We use a different notation i to represent that the cells in the scCAS data are different from the cells in the scRNA-Seq data.

x_{ig} denotes the observed gene score for gene g in cell i . The gene score summarizes the accessibility of the regions around the gene body (Cusanovich et al., 2018a). We model it by a mixture distribution with density functions $f_1(\cdot)$, $f_0(\cdot)$ and binary latent variable o_{ig} . $o_{ig} = 1$, and 0 represent the mixture components with high (f_1) and low (f_0) gene scores, respectively. Accessibility tends to be positively associated with activity of the gene. We model this positive relationship by the distribution $o_{ig}|u_{ig}$. When gene g is active in cell i ($u_{ig} = 1$), the probability that it has high gene score ($o_{ig} = 1$) is π_{i1} ; When gene g is inactive in cell i ($u_{ig} = 0$), the probability that it has high gene score ($o_{ig} = 1$) is π_{i0} . We assume that $\pi_{i1} \geq \pi_{i0}$ to represent the positive relationship. In practice, we found that fixing $\pi_{i0} = 0$ leads to good real data performance, and we set $\pi_{i0} = 0$ by default. The prior distribution $\pi_{i1} \sim \text{Beta}(\alpha = 1, \beta = 1)$. In real data example 1, the observed data is promoter accessibility and we use the same model as that for gene score.

We assume that ω_{kg}^{acc} follows Beta distribution with mean μ_{kg}^{acc} and precision ϕ^{acc} . The variable μ_{kg}^{acc} is connected with ω_{kg}^{rna} in scRNA-Seq data through the logit function: $\text{logit}(\mu_{kg}^{acc}) = f(\omega_{kg}^{rna})$. Details on the specification of $f(\cdot)$ are presented in Section 2.6.

2.3 Model for single-cell methylation data

The model specification for sc-methylation data is as the following.

$$\begin{aligned} \omega_{kg}^{met} &\xrightarrow{z_{dk}} u_{dg} \rightarrow m_{dg} \rightarrow t_{dg} \quad \forall g, \\ P(z_{dk} = 1) &= \psi_k^{met}, \\ u_{dg}|z_{dk} = 1 &\sim \text{Bernoulli}(\omega_{kg}^{met}), \\ m_{dg}|u_{dg} &\sim u_{dg}\text{Bernoulli}(\pi_{d1}) + (1 - u_{dg})\text{Bernoulli}(\pi_{d0}), \\ \pi_{d0} &\sim \text{Beta}(\alpha = 1, \beta = 1), \pi_{d1} \sim \mathbb{1}(\pi_{d1} \leq \pi_{d0})\text{Beta}(\alpha = 1, \beta = 1), \\ p(t_{dg}|m_{dg}) &= m_{dg}h_1(t_{dg}) + (1 - m_{dg})h_0(t_{dg}), \\ \omega_{kg}^{met}|\omega_{kg}^{rna} &\sim \text{Beta}(\mu_{kg}^{met}, \phi^{met}), \text{ logit}(\mu_{kg}^{met}) = g(\omega_{kg}^{rna}). \end{aligned}$$

The random variables ω_{kg}^{met} , z_{dk} , ψ_k^{met} and u_{dg} have similar interpretations to their corresponding variables in the model for scRNA-Seq data. We use a different notation d to represent that the cells in the sc-methylation data are different from the cells in the scRNA-Seq data.

The binary random variable m_{dg} denotes whether gene g is methylated in cell d , and $m_{dg} = 1$ represents that it is methylated. Methylation of a gene (promoter methylation/gene body methylation) tends to be negatively associated with activity of the gene, and we model this negative relationship with the model $m_{dg}|u_{dg}$: when

the gene g is active in cell d ($u_{dg} = 1$), it is less likely to be methylated ($m_{dg} = 1$), as we assume that $\pi_{d1} \leq \pi_{d0}$.

t_{dg} denotes the observed methylation level for gene g in cell d , and we assume that $t_{dg}|m_{dg}$ follows a mixture distribution, where $h_1(\cdot)$ and $h_0(\cdot)$ are density functions conditional on m_{dg} . The technologies/features differ for the two real data applications to be presented: promoter methylation for the gene (Pott, 2017), and gene body methylation at non-CG sites (Luo et al., 2017).

Similar to scCAS data, we connect μ_{kg}^{met} , which is the mean of ω_{kg}^{met} , and ω_{kg}^{rna} through the logit function: $\text{logit}(\mu_{kg}^{met}) = g(\omega_{kg}^{rna})$. Details on specification of $g(\cdot)$ are presented in Section 2.6.

2.4 More on model specification

Methylation and chromatin accessibility regulate gene expression biologically. Our model is specified in the reverse order, so gene expression plays a central role. This is because scRNA-Seq data is usually less noisy compared with scCAS data and sc-methylation data, the model specified this way will improve the clustering performance of scCAS data and sc-methylation data, without sacrificing much the clustering performance of scRNA-Seq data.

2.5 Prior specifications

We assume the following priors for ψ_k^{acc} , ψ_k^{rna} , ψ_k^{met} , ω_{kg}^{rna} .

$$\psi_k^{acc} \sim \text{Dir}(2, \dots, 2), \psi_k^{rna} \sim \text{Dir}(2, \dots, 2), \psi_k^{met} \sim \text{Dir}(2, \dots, 2),$$

$$\omega_{kg}^{rna} \sim \text{Beta}(\alpha_1 = 2, \beta_1 = 2)$$

The prior specification $\text{Beta}(\alpha = 2, \beta = 2)$ improves the stability of the EM algorithm in Section 3 over uniform distribution.

2.6 Determination of $f(\omega_{kg}^{rna})$, $g(\omega_{kg}^{rna})$, ϕ^{acc} and ϕ^{met}

We assume that $f(\omega_{kg}^{rna}) = \eta + \gamma\omega_{kg}^{rna} + \tau(\omega_{kg}^{rna})^2$ and $g(\omega_{kg}^{rna}) = \delta + \theta\omega_{kg}^{rna}$. The parameters $\{\eta, \gamma, \tau, \delta, \theta, \phi^{acc}, \phi^{met}\}$ are estimated empirically from the datasets. We first set the number of clusters $K=1$ and use the model to estimate ω_{kg}^{rna} , ω_{kg}^{acc} and ω_{kg}^{met} separately without considering the links on ω across the three datasets, and then fix $\hat{\omega}_{kg}^{rna}$, $\hat{\omega}_{kg}^{acc}$ and $\hat{\omega}_{kg}^{met}$ to estimate $\{\eta, \gamma, \tau, \delta, \theta, \phi^{acc}, \phi^{met}\}$ by beta regression (Silvia and Francisco, 2004). The rationale for fixing $K=1$ to estimating the parameters in the functions $f(\cdot)$ and $g(\cdot)$ that link the three modalities is that the majority of the features may not change much across the cell types. We fix $\{\hat{\eta}, \hat{\gamma}, \hat{\tau}, \hat{\delta}, \hat{\theta}, \hat{\phi}^{acc}, \hat{\phi}^{met}\}$ when implementing the EM algorithm in Section 3. Estimating $\{\eta, \gamma, \tau, \delta, \theta, \phi^{acc}, \phi^{met}\}$ separately from the EM algorithm improves computational efficiency and avoids problematic local modes. Distributions of $\hat{\omega}_{kg}^{rna}$ v.s. $\hat{\omega}_{kg}^{acc}$ and $\hat{\omega}_{kg}^{rna}$ v.s. $\hat{\omega}_{kg}^{met}$ for the two real data applications are presented in Supplementary Figures S4 and S5, we can see from Supplementary Figures S4 and S5 that the linear and quadratic models capture the trends on how $\hat{\omega}_{kg}^{acc}$ and $\hat{\omega}_{kg}^{met}$ changes with $\hat{\omega}_{kg}^{rna}$.

2.7 The mixture components

For scCAS data, we apply $f_1(x) = 0, f_0(x) = 1$ if $x=0$ and $f_1(x) = 1, f_0(x) = 0$ if $x > 0$, due to the sparsity of the data matrix.

For scRNA-Seq data, we first normalize read counts to TPM (transcripts per million) or FPKM (fragments per kilobase of exon model per million reads mapped) to account for sequencing depth and gene length, then fit a two-component gamma mixture model for the non-zero entries, through pooling $\ln(\text{TPM} + 1)$ or $\ln(\text{FPKM} + 1)$ over all the samples, and then the remaining zero entries are merged with the mixture component that has a smaller mean. The log transformation takes into account the very large values in the data matrix.

sc-methylation data represents the proportion of methylated sites within a given genomic interval, where the entries in the data matrix take values between 0 and 1. In the two real data applications, majority of entries in the data matrix take small values, for the methylation data in each cell, we first divide the entries by $(1 - \text{entries})$ to map them into $[0, \infty)$. We then normalize the entries by dividing the median of non-zero entries in each cell, and then take square of the entries to boost the signals. This transformation helps to align the three modalities and it improves the clustering results significantly (Supplementary Tables S5 and S6). Because the transformed entries represent the relative evidence of the methylation status, we input the transformed entries directly as the ratio $\frac{b_1(\cdot)}{b_0(\cdot)}$ in the EM algorithm. Histograms for the distributions of the sc-methylation data are presented in Supplementary Figure S1.

2.8 Feature selection

scRNA-Seq data is usually the least noisy data type, compared with scCAS and sc-methylation data. We use scRNA-Seq data for feature selection before implementing scAMACE. We first cluster scRNA-Seq data with SC3 and then use the cluster assignments to select top 1000 features with large mean shift across different clusters. More specifically, denote the data matrix as $\mathbf{X}_{n \times p}$ (x_{ij} denotes the observation for the i th cell and j th feature), the cluster assignments as $L_{n \times 1}$ ($l_i = k$ denotes that the i th cell belongs to the k th cluster) and total number of clusters as K . For feature j , we first calculate the difference between the mean of the cells within one cell type and the mean of cells in other cell types; the differences are represented as $D(j) = (d_{1j}, \dots, d_{Kj})$, where $d_{kj} = \text{mean}_{i:l_i=k}(x_{ij}) - \text{mean}_{i:l_i \neq k}(x_{ij})$. We take the maximum entry in $D(j)$: $m(j) = \max_k D(j)$. When $m(j)$ is large, it represents that feature j has high expression in one cluster, compared with all other clusters. Finally, we select the top 1000 features with highest values in $m(j)$.

2.9 Determination of the number of clusters K

We determine the number of clusters K for the three single-cell datasets separately before we apply scAMACE. We first run K-Means for each K and calculate the average silhouette width of observations (Kaufman and Rousseeuw, 1990). Silhouette width measures how well an observation has been classified. For each observation i , the silhouette value $s(i)$ is calculated as follows. First denote by A the cluster to which observation i has been assigned and then calculate

$a(i)$ = average Euclidean distance of i to all other objects of A .

Now consider any cluster C different from A and define

$d(i, C)$ = average Euclidean distance of i to all objects of C .
 $b(i) = \min_{C \neq A} d(i, C)$.

Then $s(i) = \frac{b(i) - a(i)}{\max(b(i), a(i))}$. When cluster A contains only a single observation, we simply set $s(i) = 0$. The average of $s(i)$ for $i = 1, 2, \dots, n$ is denoted by $\bar{s}(k)$, and it is called the average silhouette width for the entire dataset. $\bar{s}(k)$ is used for the selection of K . Higher value in $\bar{s}(k)$ indicates better clustering outcome. We select K that has the maximum average silhouette width. Details for selecting K in the two real data applications are presented in Supplementary Figures S2 and S3. When the similarity of the cell types is high, the Silhouette method may choose a smaller K than the number of cell types (Supplementary Fig. S3), and we may choose a larger K instead.

3 Statistical inference: EM algorithm

Given the observed scCAS data \mathbf{X} , scRNA-Seq data \mathbf{Y} , and sc-methylation data \mathbf{T} , we treat the latent variables $\Gamma = \{Z, U, O, V, M\}$ as missing data, and use the EM algorithm to estimate the parameters $\Phi = \{\psi^{acc}, \omega^{acc}, \pi_i, \psi^{ma}, \omega^{ma}, \pi_i, \psi^{met}, \omega^{met}, \pi_d\}$. The Q-function is $Q(\Phi | \Phi_{old}) = \mathbb{E}_{old}(\ln(P(\Phi, \Gamma | obs.)))$, where the expectation is over Γ under distribution $P(\Gamma | \Phi_{old}, obs.)$.

In the M-step, we maximize $Q(\Phi | \Phi_{old})$ with respect to Φ and update parameters as follows.

$$\begin{aligned} \psi_k^{acc} &= \frac{1 + \sum_i \mathbb{E}_{old}(z_{ik})}{K + n_{acc}}, \\ \psi_k^{ma} &= \frac{1 + \sum_i \mathbb{E}_{old}(z_{lk})}{K + n_{ma}}, \\ \psi_k^{met} &= \frac{1 + \sum_d \mathbb{E}_{old}(z_{dk})}{K + n_{met}}, \\ \pi_{i1} &= \frac{\sum_k \sum_g \mathbb{E}_{old}(z_{ik} u_{ig} o_{ig}) + \alpha_{acc} - 1}{\sum_k \sum_g \mathbb{E}_{old}(z_{ik} u_{ig}) + \alpha_{acc} + \beta_{acc} - 2}, \\ \pi_{i1} &= \frac{\sum_k \sum_g \mathbb{E}_{old}(z_{lk} u_{lg} v_{lg})}{\sum_k \sum_g \mathbb{E}_{old}(z_{lk} u_{lg})}, \\ \pi_{i0} &= \frac{\sum_k \sum_g \mathbb{E}_{old}[z_{lk}(1 - u_{lg})v_{lg}]}{\sum_k \sum_g \mathbb{E}_{old}[z_{lk}(1 - u_{lg})]}, \\ \pi_{d1} &= \frac{\sum_k \sum_g \mathbb{E}_{old}(z_{dk} u_{dg} m_{dg})}{\sum_k \sum_g \mathbb{E}_{old}(z_{dk} u_{dg})}, \\ \pi_{d0} &= \frac{\sum_k \sum_g \mathbb{E}_{old}[z_{dk}(1 - u_{dg})m_{dg}] - 1}{\sum_k \sum_g \mathbb{E}_{old}[z_{dk}(1 - u_{dg})] - 1}, \\ \omega_{kg}^{acc} &= \frac{\sum_i \mathbb{E}_{old}(z_{ik} u_{ig}) + \mu_{kg}^{acc} \phi^{acc} - 1}{\sum_i \mathbb{E}_{old}(z_{ik}) + \phi^{acc} - 2}, \\ \omega_{kg}^{met} &= \frac{\sum_d \mathbb{E}_{old}(z_{dk} u_{dg}) + \mu_{kg}^{met} \phi^{met} - 1}{\sum_d \mathbb{E}_{old}(z_{dk}) + \phi^{met} - 2}. \end{aligned}$$

We use grid search to update ω_{kg}^{ma} because its optimal value does not have an explicit form.

We iterate between E-step and M-step until converge. $\mathbb{E}(Z_i)$, $\mathbb{E}(U_i)$ and $\mathbb{E}(Z_d)$ in the last iteration are used for clustering. Details for the derivations are presented in Supplementary Materials.

4 Simulation studies

To validate scAMACE, we generated three different types of simulated data \mathbf{x} , \mathbf{y} and \mathbf{t} following the model assumption. In the simulated data, the sample sizes $n_x = 900$, $n_y = 1100$ and $n_t = 1000$. The number of features $p = 1000$. The number of clusters $K_x = K_y = K_t = 3$, and $\psi_x = \psi_y = \psi_t = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$. $f(\omega_{kg}^y) = \eta + \gamma \omega_{kg}^y + \tau (\omega_{kg}^y)^2 = -1 + 7\omega_{kg}^y - 2(\omega_{kg}^y)^2$, $g(\omega_{kg}^y) = \delta + \theta \omega_{kg}^y = -2 + 5\omega_{kg}^y$, $\phi^x = 10$ and $\phi^t = 10$. The detailed simulation scheme is presented in Supplementary Materials.

For the first data type, \mathbf{x} , we set $f_1(x) = 0$ if $x = 0$, and $f_0(x) = 0$ if $x = 1$. We fit a two-component gamma mixture model for \mathbf{y} using ‘gammamixEM’ in R (Young et al., 2019) and beta mixture model for \mathbf{t} using ‘betamix’ in R (Cribari-Neto and Zeileis, 2010; Grun et al., 2012) to estimate the mixture densities. We apply the method in Section 2.6 to estimate parameters in $f(\omega_{kg}^y)$ and $g(\omega_{kg}^y)$. We then implement scAMACE using the estimated densities and $\hat{\eta}, \hat{\gamma}, \hat{\tau}, \hat{\phi}^x, \hat{\delta}, \hat{\theta}, \hat{\phi}^t$.

We use purity, rand index, adjusted rand index and normalized mutual information to evaluate the clustering results. We implement scAMACE either on the three data types separately [‘scAMACE (separate)’] without borrowing information or jointly [‘scAMACE (joint)’]. Table 1 presents the simulation results. We also compared scAMACE with other existing methods under four additional simulation schemes: imbalanced datasets where the number of cells varies across the three datasets (Supplementary Table S1), unequal numbers of clusters in the three datasets (Supplementary Table S2), imbalanced cluster sizes (Supplementary Table S3) and smaller number of features (Supplementary Table S4). scAMACE performs the best compared with the other methods in the above simulation settings. This is likely due to integration of information from all three datasets.

In the following two real data applications, we apply methods mentioned in Section 2.7 instead of fitting a beta mixture model to sc-methylation data.

Table 1. Mean and SD (in parentheses) of purity, rand index, adjusted rand index (ARI) and normalized mutual information (NMI) for 50 independent runs are shown

	Data type	Purity	Rand index	ARI	NMI
scAMACE (joint)	<i>x</i>	0.690(0.025)	0.683(0.018)	0.288(0.041)	0.245(0.035)
	<i>y</i>	0.897(0.009)	0.874(0.010)	0.716(0.022)	0.637(0.022)
	<i>t</i>	0.704(0.021)	0.693(0.016)	0.310(0.034)	0.265(0.030)
scAMACE (separate)	<i>x</i>	0.659(0.028)	0.662(0.018)	0.241(0.041)	0.205(0.034)
	<i>y</i>	0.838(0.012)	0.810(0.012)	0.573(0.028)	0.498(0.026)
	<i>t</i>	0.643(0.020)	0.651(0.012)	0.216(0.028)	0.185(0.024)
K-means	<i>x</i>	0.383(0.020)	0.558(0.004)	0.007(0.008)	0.008(0.007)
	<i>y</i>	0.714(0.036)	0.702(0.026)	0.331(0.058)	0.283(0.049)
	<i>t</i>	0.388(0.021)	0.560(0.004)	0.010(0.008)	0.106(0.008)
Hierarchical clustering	<i>x</i>	0.360(0.011)	0.488(0.047)	0.001(0.001)	0.003(0.002)
	<i>y</i>	0.366(0.011)	0.520(0.026)	0.002(0.002)	0.003(0.002)
	<i>t</i>	0.360(0.010)	0.532(0.022)	0.001(0.001)	0.002(0.001)
Spectral clustering	<i>x</i>	0.395(0.020)	0.561(0.004)	0.012(0.009)	0.013(0.008)
	<i>y</i>	0.722(0.018)	0.708(0.018)	0.344(0.041)	0.295(0.034)
	<i>t</i>	0.400(0.025)	0.562(0.005)	0.014(0.012)	0.015(0.011)

5 Application to real data

5.1 Application 1: K562 and GM12878 scRNA-Seq, scATAC-Seq and sc-methylation data

We evaluate scAMACE by jointly clustering scRNA-Seq, scATAC-Seq and sc-methylation data generated from two cell types, K562 and GM12878 (Buenrostro *et al.*, 2015; Li *et al.*, 2017; Pott, 2017). We set $K = 2$, and use the true cell labels as a benchmark to evaluate the performance of the clustering methods. Tables 2 and 3 presents the clustering results. scAMACE performs well in separating the cell types. scRNA-Seq is perfectly separated, while there are only three cells that are not classified correctly in the sc-methylation dataset and eleven misclassifications in the scATAC-Seq dataset. In addition, the two cell types are correctly matched across the three datasets. Compared with the clustering results given by implementing scAMACE separately on the three datasets, jointly clustering the

Table 2. Clustering tables for K562, GM12878 scRNA-Seq, scATAC-Seq and sc-methylation data

		scAMACE (joint)		scAMACE (separate)				
		1	2	1	2			
scATAC-Seq	GM12878	368	5	254	119			
	K562	6	660	171	495			
scRNA-Seq	GM12878	128	0	128	0			
	K562	0	73	0	73			
sc-methyl	GM12878	16	3	7	12			
	K562	0	11	11	0			
		Seurat V3			LIGER			
		1	2	3	1	2	3	
scATAC-Seq	GM12878	346	27					
	K562	499	167					
scRNA-Seq	GM12878	101	2	25	127	0	1	
	K562	0	73	0	10	63	0	
sc-methyl	GM12878						19	
	K562						11	

Table 3. Comparison of the performance of different methods on the K562, GM12878 dataset by adjusted rand index

	scAMACE (joint)	scAMACE (separate)	Seurat V3	LIGER	scMC
scATAC-Seq	0.958	0.192	0.033		0.000
scRNA-Seq	1.000	1.000	0.713	0.800	0.771
sc-methyl	0.628	0.260		0.000	0.000

three datasets improves the overall clustering performance, especially for scATAC-Seq data, which is likely due to the integration of information across the three datasets.

We compared scAMACE with Seurat V3 (Stuart *et al.*, 2019), LIGER (Welch *et al.*, 2019) and scMC (Zhang and Nie, 2021), which are methods for integrative analysis of single-cell data. Examples were presented in Seurat V3 (Stuart *et al.*, 2019) where scRNA-Seq and scATAC-Seq data were integrated. So we implemented Seurat V3 to integrate these two data types. Seurat V3 did not perform well for scATAC-Seq data (Table 2). Seurat V3 is not applicable to integrate sc-methylation data with the other two datasets. Examples were presented in LIGER (Welch *et al.*, 2019) where scRNA-Seq data and sc-methylation data were integrated. So we implemented LIGER to integrate these two data types. LIGER did not perform well on sc-methylation data (Table 2). We also implemented LIGER to integrate all three datasets, and LIGER still did not perform well on sc-methylation data (Supplementary Tables S7 and S8), this may be due to the small sample size in sc-methylation data. scMC (Zhang and Nie, 2021) was developed for the integrative analysis of multiple single-cell datasets with the same data type. Since the features in scATAC-Seq data, scRNA-Seq data and sc-methylation data are linked, scMC can be implemented in principle. scMC did not perform well on scATAC-Seq data and sc-methylation data (Supplementary Tables S7 and S8). This may be due to the fact that the characteristics of different data types are very different, and ignoring the difference leads to suboptimal performance.

5.2 Application 2: mouse neocortex scRNA-Seq, sci-ATAC-Seq and sc-methylation data

In this example, we evaluate scAMACE for the joint analysis of single-cell datasets where the cell types are different across the datasets.

We collected single-cell datasets generated from mouse neocortex. There are five cell types in scRNA-Seq data (Tasic *et al.*, 2018), including astrocytes (Astro), glutamatergic neurons in layer 4 (L4), corticothalamic glutamatergic neurons in layer 6 (L6 CT), oligodendrocytes (Oligo) and Pvalb+ GABAergic neurons (Pvalb). There are three cell types in sci-ATAC-Seq data (Cusanovich *et al.*, 2018b), including astrocytes (Astro), excitatory neurons CPN (Ex. neurons

CPN) and oligodendrocytes (Oligo). There are three cell types in sc-methylation dataset (Luo *et al.*, 2017), including excitatory neurons in layer 4 (L4), excitatory neurons in layer 6 [labeled as L6-2 in (Luo *et al.*, 2017)] and Pvalb+ GABAergic neurons (Pvalb). In the three datasets, the optimal numbers of clusters chosen by the Silhouette method, $\hat{K} = 2$, tend to be smaller than the numbers of cell types, which is likely due to the similarity of the neuronal subtypes. We set $K=5$ when we implement scAMACE, instead of the value given by the Silhouette method. The true cell labels are used as a benchmark for evaluating the performance of the clustering methods.

The clustering results are presented in Tables 4 and 5. Even though K is larger than the number of cell types in sci-ATAC-Seq data and sc-methylation data, scAMACE still determines the correct number of cell types in sci-ATAC-Seq data. Although the cells in sc-methylation data fall into four clusters, there are only seven cells in cluster 4. Cell types in all three datasets are well separated. Astrocytes and oligodendrocytes are matched across scRNA-Seq data and sci-ATAC-Seq data. Excitatory neurons CPN in sci-ATAC-Seq data are matched with glutamatergic neurons in layer 4 in the scRNA-Seq data. We note that most excitatory neurons are glutamatergic neurons. Excitatory neurons in layers 4 and 6, and Pvalb+ GABAergic neurons are matched between scRNA-Seq data and sc-methylation data.

Compared with implementing scAMACE on the three datasets separately, the joint analysis leads to improvement in clustering, especially for sc-methylation dataset. This is likely because the joint model borrows information across the three datasets. Similar to application 1, we implemented Seurat V3 to integrate scRNA-Seq and

sci-ATAC-Seq data. Seurat V3 (Stuart *et al.*, 2019) does not perform well on sci-ATAC-Seq data (Table 4). We implemented LIGER (Welch *et al.*, 2019) to integrate scRNA-Seq and sc-methylation data. LIGER does not separate excitatory neurons in layer 4 and layer 6 in sc-methylation data (Table 4). We also integrated all three datasets by LIGER (Welch *et al.*, 2019) and scMC (Zhang and Nie, 2021). LIGER and scMC did not perform well (Supplementary Tables S9 and S10). Overall, scAMACE performed the best compared with the other methods.

5.3 Computational cost

LIGER, Seurat V3 and scMC only provide the versions that are implemented on CPU, while scAMACE can be implemented on both CPU and GPU. We summarized the computational time for scAMACE (CPU version and GPU version in python), LIGER (Welch *et al.*, 2019), Seurat V3 (Stuart *et al.*, 2019) and scMC (Zhang and Nie, 2021) (Supplementary Tables S11–S13). We implemented scAMACE, LIGER and scMC to cluster the three types of data simultaneously, and we implemented Seurat V3 to cluster scCAS data and scRNA-Seq data.

On real data application 2 (~8000 cells), the computational time for scAMACE are 418.858 s on one 3.4 GHz Intel Xeon Gold CPU and 69.652 s on one 3.1 GHz Dual Intel Xeon Gold GPU. Compared with LIGER (80.389 s on one 3.4 GHz Intel Xeon Gold CPU), scMC (372.323 s on one 3.4 GHz Intel Xeon Gold CPU) and Seurat V3 (116.688 s for scRNA-Seq and sci-ATAC-Seq data on one

Table 4. Clustering tables for the mouse neocortex scRNA-Seq, sci-ATAC-Seq and sc-methylation data

		scAMACE (joint)					scAMACE (separate)										
		1	2	3	4	5	1	2	3	4	5						
sci-ATAC-Seq	Astro	550	0		1		550	0	1								
	Ex. neurons CPN	0	1391		0		1	1390	0								
	Oligo	0	1		457		0	0	458								
scRNA-Seq	Astro	368	0	0	0	0	368	0	0	0	0						
	L4	0	1401	0	0	0	0	1401	0	0	0						
	L6 CT	0	0	960	0	0	0	0	960	0	0						
	Oligo	25	0	0	66	0	27	0	0	64	0						
	Pvalb	0	0	0	0	1337	0	0	0	0	1337						
sc-methyl	L4		411	1	0	0	412										
	L6-2		20	703	6	0	729										
	Pvalb		0	0	1	153	154										
		Seurat V3									LIGER						
		1	2	3	4	5	6	7	8	9	1	2	3	4	5	6	7
sci-ATAC-Seq	Astro	296	29	1	2	223											
	Ex. neurons CPN	110	461	243	0	577											
	Oligo	64	90	1	0	303											
scRNA-Seq	Astro	0	0	0	0	0	368	0	0	0	6	0		0	362		0
	L4	1028	0	0	0	373	0	0	0	0	0	1401		0	0		0
	L6 CT	0	960	0	0	0	0	0	0	0	2	0		958	0		0
	Oligo	0	0	0	0	0	0	0	60	31	31	0		0	0		60
	Pvalb	0	0	647	498	0	0	192	0	0	1337	0		0	0		0
sc-methyl	L4										0	11	679				0
	L6-2										0	13	399				0
	Pvalb										68	0	0				86

Table 5. Comparison of the performance of different methods on the mouse neocortex dataset by adjusted rand index

	scAMACE (joint)	scAMACE (separate)	Seurat V3	LIGER	scMC
sci-ATAC-Seq	0.998	0.998	0.058		0.019
scRNA-Seq	0.997	0.997	0.697	0.983	0.145
sc-methyl	0.932	0.000		0.316	0.001

3.4 GHz Intel Xeon Gold CPU), scAMACE has competitive computational speed, especially the GPU version.

Next, we generated a dataset with sample size = 30,000 ($n_{acc} = n_{ma} = n_{met} = 10,000$) by sampling the cells with replacement from real data application 2. The computational time for scAMACE are 1534.631 s on one 3.4 GHz Intel Xeon Gold CPU and 250.089 s on one 3.1 GHz Dual Intel Xeon Gold GPU. Compared with LIGER (555.574 s on one 3.4 GHz Intel Xeon Gold CPU), scMC (3667.878 s on one 3.4 GHz Intel Xeon Gold CPU) and Seurat V3 (290.640 s for scRNA-Seq and sci-ATAC-Seq data on one 3.4 GHz Intel Xeon Gold CPU), scAMACE has competitive computational speed on datasets with larger scale.

6 Conclusion

Unsupervised methods including dimension reduction and clustering are essential to the analysis of single-cell genomic data as the cell types are usually unknown. We have developed scAMACE, a model-based approach for integratively clustering single-cell data on chromatin accessibility, gene expression and methylation. scAMACE provides statistical inference of cluster assignments and achieves better cell type separation combining biological information across different types of genomic features. In the two real data applications, the scRNA-Seq data are generated from the SMART-Seq platform (Li et al., 2017; Tasic et al., 2018). To implement scAMACE on UMI-based scRNA-Seq data (10x data), we may need to modify the distributions of the mixture components $g_0(\cdot)$ and $g_1(\cdot)$. The cells in our real data examples are differentiated and mature cells. In the future, we will investigate the performance of scAMACE on immature cells undergoing differentiation.

Data availability

Real data application 1 human K562-GM12878 scRNA-Seq data was retrieved from NCBI Gene Expression Omnibus (GEO) with the accession number GSE81861, scATAC-Seq data was collected from GEO with the accession number GSE65360, and sc-methylation data is available at GEO with the accession number GSE83882. Real data application 2 mouse neocortex scRNA-Seq data are available at GEO with the accession number GSE115746, scATAC-Seq datasets were downloaded from <https://atlas.gs.washington.edu/mouse-atac/data/>, and sc-methylation data are available at GEO with the accession number GSE97179.

Acknowledgements

The authors thank the reviewers' constructive comments and suggestions, and Jinwen Yang, Wenyu Zhang, Pengcheng Zeng for the helpful discussions.

Funding

This work was supported by the Chinese University of Hong Kong direct grants [4053360, 4053423], the Chinese University of Hong Kong startup grant [4930181], the Chinese University of Hong Kong's Project Impact Enhancement Fund (PIEF) and Science Faculty's Collaborative Research Impact Matching Scheme (CRIMS) and Hong Kong Research Grant Council [ECS 24301419, GRF 14301120].

Conflict of interest: The authors declare that they have no conflict of interest.

References

Bravo González-Blas, C. et al. (2019) cisTopic: cis-regulatory topic modeling on single-cell ATAC-seq data. *Nat. Methods*, **16**, 397–400.
 Buenrostro, J.D. et al. (2015) Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature*, **523**, 486–490.
 Corces, M.R. et al. (2016) Lineage-specific and single-cell chromatin accessibility charts human hematopoiesis and leukemia evolution. *Nat. Genet.*, **48**, 1193–1203.
 Cribari-Neto, F. and Zeileis, A. (2010) Beta regression in R. *J. Stat. Softw.*, **34**, 1–24.

Cusanovich, D.A. et al. (2018a) The cis-regulatory dynamics of embryonic development at single-cell resolution. *Nature*, **555**, 538–542.
 Cusanovich, D.A. et al. (2018b) A single-cell atlas of in vivo mammalian chromatin accessibility. *Cell*, **174**, 1309–1324.e18.
 Duren, Z. et al. (2017) Modeling gene regulation from paired expression and chromatin accessibility data. *Proc. Natl. Acad. Sci. USA*, **114**, E4914–E4923.
 Duren, Z. et al. (2018) Integrative analysis of single-cell genomics data by coupled nonnegative matrix factorizations. *Proc. Natl. Acad. Sci. USA*, **115**, 7723–7728.
 Grun, B. et al. (2012) Extended beta regression in R: shaken, stirred, mixed, and partitioned. *J. Stat. Softw.*, **48**, 1–25.
 Hui, T. et al. (2018) High-resolution single-cell DNA methylation measurements reveal epigenetically distinct hematopoietic stem cell subpopulations. *Stem Cell Rep.*, **11**, 578–592.
 Kapourani, C.-A. and Sanguinetti, G. (2016) Higher order methylation features for clustering and prediction in epigenomic studies. *Bioinformatics*, **32**, i405–i412.
 Kapourani, C.-A. and Sanguinetti, G. (2019) Melissa: Bayesian clustering and imputation of single-cell methylomes. *Genome Biol.*, **20**, 61.
 Kaufman, L. and Rousseeuw, P.J. (1990) *Finding Groups in Data: An Introduction to Cluster Analysis*. John Wiley, New York.
 Kiselev, V.Y. et al. (2017) Sc3: consensus clustering of single-cell RNA-seq data. *Nat. Methods*, **14**, 483–486.
 Lahmehmann, D. et al. (2020) Eleven grand challenges in single-cell data science. *Genome Biol.*, **21**, 31.
 Li, H. et al. (2017) Reference component analysis of single-cell transcriptomes elucidates cellular heterogeneity in human colorectal tumors. *Nat. Genet.*, **49**, 708–718.
 Lin, P. et al. (2017) CIDR: ultrafast and accurate clustering through imputation for single-cell RNA-seq data. *Genome Biol.*, **18**, 59.
 Lin, Z. et al. (2020) Model-based approach to the joint analysis of single-cell data on chromatin accessibility and gene expression. *Statist. Sci.*, **35**, 2–13.
 Luo, C. et al. (2017) Single-cell methylomes identify neuronal subtypes and regulatory elements in mammalian cortex. *Science*, **357**, 600–604.
 Ma, A. et al. (2020) Integrative methods and practical challenges for single-cell multi-omics. *Trends Biotechnol.*, **38**, 1007–1022.
 Ng, A.Y. et al. (2001). On spectral clustering: analysis and an algorithm. In *Proceedings of the 14th International Conference on Neural Information Processing Systems: Natural and Synthetic, NIPS'01*. MIT Press, Cambridge, MA, USA, pp. 849–856.
 Pott, S. (2017) Simultaneous measurement of chromatin accessibility, DNA methylation, and nucleosome phasing in single cells. *eLife*, **6**, e23203.
 Silvia, F. and Francisco, C.-N. (2004) Beta regression for modelling rates and proportions. *J. Appl. Stat.*, **31**, 799–815.
 Stuart, T. et al. (2019) Comprehensive integration of single-cell data. *Cell*, **177**, 1888–1902.e21.
 Tasic, B. et al. (2018) Shared and distinct transcriptomic cell types across neocortical areas. *Nature*, **563**, 72–78.
 Wang, B. et al. (2017) Visualization and analysis of single-cell RNA-seq data by kernel-based similarity learning. *Nat. Methods*, **14**, 414–416.
 Wang, C. et al. (2020) Integrative analyses of single-cell transcriptome and regulome using maestro. *Genome Biol.*, **21**, 198.
 Ward, J.H. (1963) Hierarchical grouping to optimize an objective function. *J. Am. Stat. Assoc.*, **58**, 236–244.
 Welch, J.D. et al. (2019) Single-cell multi-omic integration compares and contrasts features of brain cell identity. *Cell*, **177**, 1873–1887.e17.
 Xiong, L. et al. (2019) Scale method for single-cell atac-seq analysis via latent feature extraction. *Nat. Commun.*, **10**, 4576.
 Yin, L. et al. (2019) Virtual methylome dissection facilitated by single-cell analyses. *Epigenet. Chromatin*, **12**, 66.
 Young, D.S. et al. (2019) Finite mixture-of-gamma distributions: estimation, inference, and model-based clustering. *Adv. Data Anal. Classif.*, **13**, 1053–1082.
 Zamanighomi, M. et al. (2018) Unsupervised clustering and epigenetic classification of single cells. *Nat. Commun.*, **9**, 2410.
 Zeng, P. et al. (2020) Coupled co-clustering-based unsupervised transfer learning for the integrative analysis of single-cell genomic data. *Brief. Bioinform.*, **bbaa347**.
 Zhang, L. and Nie, Q. (2021) SCMC learns biological variation through the alignment of multiple single-cell genomics datasets. *Genome Biol.*, **22**, 10.
 Zhu, L. et al. (2019) Semisoft clustering of single-cell data. *Proc. Natl. Acad. Sci. USA*, **116**, 466–471.